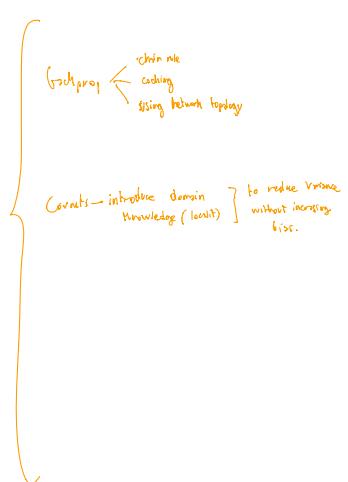


→ Problematis → Model → Optimization Problem →
 → Vectorizing to simplify math / speed computation
 → Poly features
 → Kernelization
 → Bias - Variance Tradeoff



Non-linear Least Squares

Until now just considering predictions for $f(x)$ that are linear in x . However a lot of problems would require non-linear models in their natural form. We could do feature extraction, but it's not equivalent.

Let's consider an arbitrary model $\{f(x)\}$ (nonlinear) that is not necessarily linear. Note that the likelihood of the data will be:

$$\text{Assume } Y_i \sim \mathcal{N}(f(x_i), \sigma^2) \Rightarrow Y_i | x_i \sim \mathcal{N}(f(x_i), \sigma^2)$$

$$\begin{aligned} \hat{\theta}_{\text{OLS}} &= \underset{\theta}{\operatorname{argmin}} \sum_i \log(Y_i | x_i) = \underset{\theta}{\operatorname{argmin}} \sum_i \log(f(x_i)) \\ &= \underset{\theta}{\operatorname{argmin}} \sum_i \left(Y_i - f(x_i) \right)^2 = \underset{\theta}{\operatorname{argmin}} E(\theta) \end{aligned}$$

\Rightarrow By first order optimality conditions we know that global minimum must occur at a point where gradient is zero (no constraint)

$$\begin{aligned} \Rightarrow \nabla_{\theta} E(\theta) &= \frac{\partial}{\partial \theta} \left(Y_i - f(x_i) \right)^2 = 0 \\ &= (Y - f(X))^T (Y - f(\theta)) = 0 \quad (\text{in more compact notation}) \\ &= J(\theta)^T (Y - F(\theta)) = 0 \quad J(\theta) = \begin{bmatrix} f(\theta)^T \\ \vdots \\ f(\theta)^T \end{bmatrix}, F(\theta) = \begin{bmatrix} f(x_1) \\ \vdots \\ f(x_n) \end{bmatrix} \quad \text{leads to OLS solution} \end{aligned}$$

If not conv, no closed solution (multiple points)
Even when convex, not always closed

Gauss-Newton Algorithm

Must resort to an iterative algorithm \rightarrow GN at each step linearly approximates $F(\theta)$ for current value of θ

Here's guess θ at iteration K , θ^K
ensures that J can only depend on x_i , so reduces the problem that we can solve for θ (if we want to) one

$$\begin{aligned} \rightarrow F(\theta) &\approx \tilde{F}(\theta) = F(\theta^K) + \frac{\partial}{\partial \theta} F(\theta^K)(\theta - \theta^K) \\ &= F(\theta^K) + J(\theta^K)(\theta - \theta^K) \end{aligned}$$

\rightarrow Now by FOC as above we get

$$\begin{aligned} J(\theta^K)^T (Y - \tilde{F}(\theta)) &= J(\theta^K)^T [Y - F(\theta^K) - J(\theta^K) \Delta \theta] \\ &\quad \text{or} \\ &= J^T [\Delta Y - J \Delta \theta] = 0 \Rightarrow J^T \Delta Y = J^T J \Delta \theta \\ &\Rightarrow \theta^K = (\mathbf{J}^T \mathbf{J})^{-1} \mathbf{J}^T \Delta Y \\ &\Rightarrow \theta^* = \theta^K + (\mathbf{J}^T \mathbf{J})^{-1} \mathbf{J}^T \Delta Y \quad \checkmark \text{optional: B guess the approx. F} \end{aligned}$$

$$\therefore \theta^{*m} = \theta^* + \Delta \theta = \theta^*$$

Optimization

The general optimization problem will be of the form $\min_{x \in X} f(x)$
 $\left. \begin{array}{l} \text{most optimization does not have} \\ \text{closed form solutions} \end{array} \right\} \text{we just use iterative algorithms}$

Strong Convexity

A function is **strongly convex** iff any of these equivalent conditions hold:

- $\{f(x + t\beta) \leq f(x) + t\langle \nabla f(x), \beta \rangle - \frac{\kappa t^2}{2} \|y - x\|^2 \quad \forall y, \beta, t \geq 0\}$
- $\{f(x + t\beta) \leq f(x) + t\langle \nabla f(x), \beta \rangle - \frac{\kappa t^2}{2} \|\nabla f(x)\|^2 \quad \forall y, \beta, t \geq 0\}$

$$V: \nabla^2 f(x) \geq m$$

Smoothness

M -smooth \Rightarrow Lipschitz continuous gradient function iff:

$$\|\nabla f(y) - \nabla f(x)\| \leq M \|y - x\|$$

Gradient Descent

Theoretically takes steps in the direction of the steepest descent.

Initialize $w^{(0)}$
while $\|f(w^{(t)})\|$ not converged:
 $\quad w^{(t+1)} \leftarrow w^{(t)} - \alpha_t \nabla f(w^{(t)})$

I.4.2 Stochastic Gradient Descent	I.4.3 Coordinate Descent
Algorithm 3: Stochastic Gradient Descent Initialization: $w^{(0)}$ and learning rate α while $\ f(w^{(t)})\ $ not converged do $w^{(t+1)} \leftarrow w^{(t)} - \alpha \nabla f(w^{(t)})$ $w^{(t+1)} \leftarrow w^{(t+1)} - \alpha_t \nabla f(w^{(t)})$ \downarrow stochastic gradient descent!	Algorithm 4: Coordinate Descent Initialization: $w^{(0)}$ and learning rate α while $\ f(w^{(t)})\ $ not converged do $w^{(t+1)} \leftarrow w^{(t)} - \alpha \nabla f(w^{(t)})$ \downarrow coordinate descent!
I.4.4 Gradient Descent with Momentum	I.4.5 Conjugate Gradient Descent
Algorithm 5: Gradient Descent with Momentum Initialization: $w^{(0)}$ and learning rate α while $\ f(w^{(t)})\ $ not converged do $w^{(t+1)} \leftarrow w^{(t)} - \alpha \nabla f(w^{(t)})$ $w^{(t+1)} \leftarrow w^{(t+1)} - \alpha_t \nabla f(w^{(t)})$ \downarrow path in all directions!	Algorithm 6: Conjugate Gradient Descent Initialization: $w^{(0)}$ and learning rate α while $\ f(w^{(t)})\ $ not converged do $w^{(t+1)} \leftarrow w^{(t)} - \alpha \nabla f(w^{(t)})$ \downarrow fastest convergence!

Newton's Method

Second order approximation

Q: In NN, what do we do for 100s of training examples?
A: We will have very complex function mapping. So let's just give you a sketch from above.
We need $J(\theta)$ (gradient) \Rightarrow need to compute $J(\theta)$, $J(\theta)$, $J(\theta)$, ...
But there are many local minima. Therefore we have to use gradient descent.
We locally approximate our NN weight function $J(\theta)$ with one from a all right, and then use closed form solution. $\nabla J(\theta) = 0$
 \downarrow only changing those. Iterating complex and update until loss not over

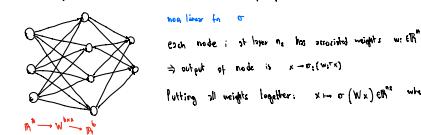
In NN, do we have to use the whole $J(\theta)$ or just gradients??
We just directly work with $J(\theta)$, i.e. $J(\theta)$ and its derivative!
What is the point of doing $J(\theta)$ if from value perspective?

Backpropagation \rightarrow gives us gradients $\frac{\partial L}{\partial \theta_i}$ \rightarrow we can then do θ_i .

$\left[\begin{array}{c} (x_1, y_1) \\ \vdots \\ (x_n, y_n) \end{array} \right]$
 \downarrow has inputs on x_i in form gradient!
we can find point for

Gradient Descent / Neural Networks

We only consider **feedforward** (DAGs). Most basic is the multilayer perceptron



↳ Idea: we always want to "fit" our model, i.e. find its parameters (for example) so that we can make predictions

non linear in x
each node: it takes x_i as its weighted inputs $w_i x_i^{(t+1)}$
 \Rightarrow output of node is $x_i - \sigma(w_i x_i^{(t+1)})$

Putting all weights together: $x \mapsto \sigma(w x) e^{w x}$ where $W = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}$

Nonlinearities
Usually σ is $tanh$, σ , softmax
However, one long exception
Sigmoid: $\theta^* = \frac{e^{\theta}}{1 + e^{\theta}}$

Therefore, we can consider the entire network to be a function:

$$x \mapsto \sigma_1(W_1) \sigma_2(\dots \sigma_n(W_n \sigma_1(W_1 x))) \rightarrow \text{neural networks are "compositionals"}$$

Expressive Power

Important to note that all the power of neural nets comes from the non-linearities without them, we would get linear grad function \rightarrow a glorified regression!

$$x \mapsto W_1 W_2 \dots W_n W_n x = W x$$

$$\text{Also } \text{rank}(W) \leq \min(\text{rank}(W_1) \dots \text{rank}(W_n))$$

Choosing non-linearity

Step function is intuitively right, but it doesn't enable us to update our weights \times gradient is always zero.

$$\text{Trying ReLU: } \sigma(x) = \begin{cases} x & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases}$$

Note that with weights a, b \times $\sigma(a + b x)$, the ReLU can shift and assume arbitrary slopes!

$$\begin{aligned} \frac{\partial L}{\partial a} &= \sum_i -2(y_i - h(x_i)) \frac{\partial h(x_i)}{\partial a} = \sum_i -2(y_i - h(x_i)) \cdot \begin{cases} 0 & \text{if } a + b x_i < 0 \\ 1 & \text{if } a + b x_i > 0 \end{cases} \\ \frac{\partial L}{\partial b} &= \sum_i -2(y_i - h(x_i)) \frac{\partial h(x_i)}{\partial b} = \sum_i -2(y_i - h(x_i)) \cdot \begin{cases} 0 & \text{if } a + b x_i < 0 \\ x_i & \text{if } a + b x_i > 0 \end{cases} \end{aligned}$$

Biased
Biases

Universal Function Approximation

If $\sigma: \mathbb{R} \rightarrow \mathbb{R}$ non-constant, bounded, non-decreasing, and continuous, and let $S \subseteq \mathbb{R}^d$ be closed and bounded. Then for any continuous $f: S \rightarrow \mathbb{R}$ and any size N NN with d hidden layer with finitely many nodes, which we can write

$$h(x) = \sum_{j=1}^N c_j \sigma(a_j + b_j x)$$

$$\therefore |h(x) - f(x)| < \epsilon \quad \forall x \in S$$

Computational Graphs

We assume our network is a DAG. This graph is called **Computational graph**.

Key: each vertex v_j is the result of a differentiable computation (can be multiple at同一时间点)

Vertex v_i is the loss function, also weights are considered as vertices

Goal: efficiently calculate $\frac{\partial L}{\partial w_i}$ $\forall i \in W$

The primary mathematical tool employed in backpropagation is chain rule:

$$\frac{\partial L}{\partial w_i} = \sum \frac{\partial L}{\partial v_j} \frac{\partial v_j}{\partial w_i} = \sum_{v_j \in \text{parents}(w_i)} \sum_{v_k \in \text{children}(v_j)} \frac{\partial L}{\partial v_k} \frac{\partial v_k}{\partial v_j} \frac{\partial v_j}{\partial w_i}$$

Main idea is to reuse components: many are the same!

Ideas: use dynamic programming to compute the derivatives following the DAG dependencies:
the main problem is computing $\partial L(w)$ and the subproblem is $\frac{\partial L}{\partial w_i}$

Examples

Fully connected layers

$$\begin{aligned} \text{Each vertex calculates: } & g_3 = \sum_i w_{ij} z_i \\ \therefore \frac{\partial L}{\partial w_{ij}} &= g_3 \\ \Rightarrow \frac{\partial L}{\partial w_{ij}} &= \frac{\partial g_3}{\partial w_{ij}} \\ \text{Also } \frac{\partial g_3}{\partial w_{ij}} &= \frac{\partial g_3}{\partial z_i} \cdot \frac{\partial z_i}{\partial w_{ij}} \end{aligned}$$

Element-wise nonlinearities

$$\begin{aligned} \text{Let } & g_i = \sigma(z_i) \\ \Rightarrow \frac{\partial L}{\partial w_{ij}} &= \frac{\partial L}{\partial g_i} \cdot \frac{\partial g_i}{\partial z_i} \cdot \frac{\partial z_i}{\partial w_{ij}} \end{aligned}$$

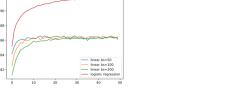


Figure 1: Epoch vs Accuracy
Figure 1 shows the plots. The training loss for 50 epochs on my laptop is 11%, 2%, and 1% respectively. The larger the batch size, the less iterations are needed to reach a certain level of accuracy. It is because for larger batch size, there is less overhead due to the vectorization. On the other hand, smaller batch size requires more iterations to reach a certain level of accuracy. So there is a trade-off here.

Implementation of backpropagation in the MNIST dataset to classify 10 digits. Let x be the feature vector and y the target vector (one-hot encoding of class, i.e., $y_i = 1$ if and only if the image is in class i and $y_i = 0$ if not).
 $\text{forward pass: } z = \sigma(w x + b)$
 $\text{backward pass: } \frac{\partial L}{\partial w} = \sum_i \frac{\partial L}{\partial z_i} \frac{\partial z_i}{\partial w} = \sum_i \frac{\partial L}{\partial z_i} \sigma'(w x + b) x$
 $\text{backward pass: } \frac{\partial L}{\partial b} = \sum_i \frac{\partial L}{\partial z_i} \frac{\partial z_i}{\partial b} = \sum_i \frac{\partial L}{\partial z_i} \sigma'(w x + b)$

Implementation of backpropagation in the MNIST dataset to classify 10 digits. Let x be the feature vector and y the target vector (one-hot encoding of class, i.e., $y_i = 1$ if and only if the image is in class i and $y_i = 0$ if not).
 $\text{forward pass: } z = \sigma(w x + b)$
 $\text{backward pass: } \frac{\partial L}{\partial w} = \sum_i \frac{\partial L}{\partial z_i} \frac{\partial z_i}{\partial w} = \sum_i \frac{\partial L}{\partial z_i} \sigma'(w x + b) x$
 $\text{backward pass: } \frac{\partial L}{\partial b} = \sum_i \frac{\partial L}{\partial z_i} \frac{\partial z_i}{\partial b} = \sum_i \frac{\partial L}{\partial z_i} \sigma'(w x + b)$

Generative vs Discriminative Classification

When dealing with the problem of classifying points in discrete categories, there are 2 main approaches:

- Generative:

$$f_K(x) = P(x|Y=k) \quad \text{distri. for class } k \\ P(Y=k) \quad \text{prior prob. for each class}$$

total of K distri. (prior + K for class)

each one $\frac{1}{K} \text{ if } x_i \in \mathcal{X}_k^i$

The idea is to assign the label with highest likelihood:

$$\hat{P}_{\text{pred}} = \arg\max_k P(x|k) = \arg\max_k \frac{\int_K(x)P(Y=k)}{\int_{\mathcal{X}}(x)} = \arg\max_k \frac{\int_K(x)P(Y=k)}{\sum_k \int_K(x)P(Y=k)} = \arg\max_k f_K(x)P(Y=k)$$

we can do this but don't read it here

Pros:
- Can generate new samples
- Data augmentation
- Quick to train

Cons:
- Inefficiency
- Num terms: $\frac{d}{2}(d+1)/2! = O(d^2)$ for our matrix
- $O(Kd^2)$ if K classes

- Discriminative:

model only $Y|X$ or not even → just have a boundary of fn f s.t. $f(X) = Y$ estimated

$\hat{y}(x)$ passes!

Least Squares SVM

An initial intuitive approach to classify points in a discriminative setting would be to try and fit a normal linear regression line between our classes (as a boundary) by setting $Y=-1, Y=1$.

Instead of considering wx as our prediction, we can consider

$$\text{sign}(y_i w^T x_i) = \begin{cases} 1 & \text{if } \text{sign}(w^T x_i) = y_i \\ -1 & \text{if } \text{sign}(w^T x_i) \neq y_i \end{cases} \Rightarrow \begin{cases} 1 & \text{if } w^T x_i > 0 \\ -1 & \text{if } w^T x_i < 0 \end{cases}$$

What does this mean geometrically? It means that we are trying to find a hyperplane (defined by w) such that all points with $y_i = 1$ lie on one side and vice versa.



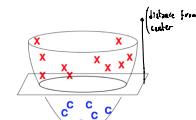
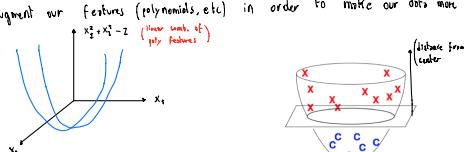
The corresponding loss function is: $\underset{w \in \mathbb{R}^{d+1}}{\text{argmin}} \sum_i (y_i - \text{sign}(w^T x_i))^2 + \lambda \|w\|_2^2$ → NP-Hard!

simplified (approximated)
 $\underset{w \in \mathbb{R}^{d+1}}{\text{argmin}} \sum_i (y_i - w^T x_i)^2 + \lambda \|w\|_2^2$
 note we are penalizing distance from what is a weird way!

Feature extension

To get non-linear decision boundaries, we can augment our features (polynomials, etc) in order to make our data more

separable:



In this case, our loss becomes:

$$\underset{w \in \mathbb{R}^{d+1}}{\text{argmin}} \sum_i (y_i - \phi(x_i)^T w_i)^2 + \lambda \|w\|_2^2$$

$$\underset{w \in \mathbb{R}^{d+1}}{\text{argmin}} \sum_i (y_i - g_w(x_i))^2 + \lambda \|w\|_2^2$$

for classification: $\hat{y}_i = \begin{cases} 1 & \text{if } g_w(x_i) > 0 \\ -1 & \text{if } g_w(x_i) \leq 0 \end{cases}$ for a some threshold

the power of NN is that we can have derivatives through this even though non-linear → very complex

Multiclass Extension

One could use the same framework for multiclass classification by using one-hot encoding. In this case we would have $y \in \mathbb{R}^K$ and $W \in \mathbb{R}^{K \times d}$

$$\text{Loss fn: } \underset{W \in \mathbb{R}^{K \times d}}{\text{argmin}} \sum_i \|y_i - Wx_i\|_2^2 + \lambda \|W\|_F^2$$

note that we are penalizing my "score" for class that are not right

One could use the same framework for multiclass classification by using one-hot encoding. In this case we would have $y \in \mathbb{R}^K$ and $W \in \mathbb{R}^{K \times d}$ i.e. $Wx \in \mathbb{R}^K$ represents our "score" for each class

Logistic Regression has no known closed form solution to minimize the loss. Therefore, one must minimize it by using SD or other stochastic forms.

Logistic Regression

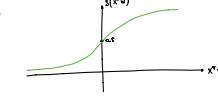
Discriminative classification method that assigns probabilities for each class (has a probabilistic interpretation, unlike SVM)

In logistic classification we use $0/1$ labels.

We use the score $x^T w$ and transform it to a probability with the sigmoid (link function) $f(z) = \frac{1}{1+e^{-z}}$

We let our prediction be

$$\hat{P} = \arg\max_k P(\hat{y}=k|x) = \begin{cases} 1 & \text{if } f(x^T w) > 0.5 \Leftrightarrow x^T w > 0 \\ 0 & \text{else} \end{cases}$$



As a loss function, we will choose cross-entropy, as

$$L(w) = -\sum_i y_i \ln \hat{P}_i + (1-y_i) \ln (1-\hat{P}_i) = -\ln P(y_i)$$

only way to have zero loss is $y_i = \hat{P}_i$

Consider the classes as coming from a Bernoulli: $P(Y_i=k) \sim \text{Bern}(p_i)$

$$P(Y_i=y_i) = p_i^{y_i} (1-p_i)^{1-y_i} \quad \text{where } p_i = \frac{1}{1+e^{-x^T w}}$$

$$\Rightarrow \hat{W}_{LR} = \arg\max_w P(Y_1=y_1, \dots, Y_n=y_n|x_1, \dots, x_n)$$

$$= \arg\min_w L(w) \rightarrow \text{Cross Entropy}$$

Information Theory Justification

KL-Divergence = $D_{KL}(P||Q) = \sum_i P(i) \ln \frac{P(i)}{Q(i)}$
 note a true distance (not symmetric) always positive (0 if $P=Q$)

$$\hat{W}_{LR} = \arg\min_w D_{KL}(P||\hat{P})$$

$$= \arg\min_w \sum_i y_i \ln \frac{y_i}{\hat{P}_i} + (1-y_i) \ln \frac{(1-y_i)}{\hat{P}_i}$$

$$= \arg\min_w -\sum_i y_i \ln \hat{P}_i + (1-y_i) \ln (1-\hat{P}_i) \quad \text{Summing the "cross-entropies" } H(\hat{P}_i, y_i)$$

$$D_{KL}(Y||\hat{Y}) = H[Y, \hat{Y}] - H(Y)$$

entropy of Y

Cross-entropy between 2 dists is # bits necessary to identify that coming from Q rather than "true" distribution P . What are Q and P in this case? why not symmetric?

$$\rightarrow H(p, q) = -\sum_x p(x) \ln q(x)$$

Multiclass Logistic Regression

In multiclass logistic we model the distribution of labels for each point to be a multinomial

$$\underset{W \in \mathbb{R}^{K \times d}}{\text{argmin}} \sum_i \text{Multi}(y_i, \dots, y_i) \quad \text{softmax fn } \hat{P}(x) = \frac{e^{w_i^T x}}{\sum_j e^{w_j^T x}}$$

where $\hat{P}_i = \frac{1}{\sum_j} \text{softmax}(Wx)$ (for true y_i , weight)

$\hat{P}_i = \sigma(w_i^T x) \quad \text{for each } j \text{ weight}$

$\hat{P}_i = \sigma(w_i^T x) \quad \dots \quad \sigma(w_K^T x)$ softmax distribution (W, x)

and $\hat{Y} = \arg\max_k \sigma(Wx)$ (probabilities sum up to one)

can the softmax be thought of as an operation that produces probabilities for our multiclass outputs

$$\text{logistic} \subseteq \text{softmax}: \begin{cases} \text{the logistic distribution is just a special case (zero) of the softmax one} \\ \text{softmax} \end{cases}$$

$$P(\hat{Y}_i=1|x_i, W) = \sigma(Wx_i)_1 = \frac{e^{w_1^T x_i}}{e^{w_1^T x_i} + e^{w_2^T x_i}} = \frac{1}{1+e^{-(w_1-w_2)^T x_i}}$$

$$P(\hat{Y}_i=0|x_i, W) = \sigma(Wx_i)_0 = \frac{e^{w_0^T x_i}}{e^{w_0^T x_i} + e^{w_1^T x_i}} = 1 - \sigma((w_1-w_0)^T x_i)$$

The loss function can be justified as before with MLE / KL divergence (as reducing sum of cross-entropy between all pairs of predicted vs true distributions)

$$L(W) = -\sum_{i=1}^n \sum_{j=1}^K \delta_{j,y_i} \log P(\hat{y}_j = j | x_i, W)$$

$$= \sum_{i=1}^n H(Y_i, \hat{Y}_i)$$

Information Approach: We have an estimated and a "true" distribution for each training point and we want to minimize the sum of the cross entropies between them

Training

Logistic regression has no known closed form solution to minimize the loss. Therefore, one must minimize it by using SD or other stochastic forms.

GD stuff.

HW 9
Last Problem

Newton's Method
Disc. g

$$\text{Assume our data comes from two classes and the prior for class 0 is } p = \pi_0. \text{ Also the prior for class 1 is } 1 - \pi_0. \text{ The log-likelihood function is } Q(w) = \ln(\pi_0) + \ln(1-\pi_0) + \sum_{i=1}^n \left[\ln P(\hat{y}_i = y_i | x_i, w) + \ln(1-P(\hat{y}_i = y_i | x_i, w)) \right]$$

$$= \left(\ln \frac{\pi_0}{1-\pi_0} \right) + \sum_{i=1}^n \left[\ln \frac{P(\hat{y}_i = y_i | x_i, w)}{1-P(\hat{y}_i = y_i | x_i, w)} \right] = -Q(w)$$

Note the first order derivative of $Q(w)$ is zero:

$$Q'(w) = \sum_{i=1}^n \frac{P(\hat{y}_i = y_i | x_i, w)}{1-P(\hat{y}_i = y_i | x_i, w)} \cdot \ln \frac{P(\hat{y}_i = y_i | x_i, w)}{1-P(\hat{y}_i = y_i | x_i, w)} = \sum_{i=1}^n \frac{P(\hat{y}_i = y_i | x_i, w)}{1-P(\hat{y}_i = y_i | x_i, w)} \cdot \frac{P(\hat{y}_i = y_i | x_i, w)}{P(\hat{y}_i = y_i | x_i, w) - P(\hat{y}_i = y_i | x_i, w)} = \sum_{i=1}^n \frac{P(\hat{y}_i = y_i | x_i, w)}{P(\hat{y}_i = y_i | x_i, w) - P(\hat{y}_i = y_i | x_i, w)} = \sum_{i=1}^n \frac{P(\hat{y}_i = y_i | x_i, w)}{P(\hat{y}_i = y_i | x_i, w)} = \sum_{i=1}^n 1 = n$$

Newton's Method: we can write it out as

$$Q''(w) = \sum_{i=1}^n \frac{P(\hat{y}_i = y_i | x_i, w)}{(1-P(\hat{y}_i = y_i | x_i, w))^2} \cdot \ln \frac{P(\hat{y}_i = y_i | x_i, w)}{1-P(\hat{y}_i = y_i | x_i, w)} + \sum_{i=1}^n \frac{P(\hat{y}_i = y_i | x_i, w)}{(1-P(\hat{y}_i = y_i | x_i, w))^2} \cdot \frac{P(\hat{y}_i = y_i | x_i, w)}{P(\hat{y}_i = y_i | x_i, w) - P(\hat{y}_i = y_i | x_i, w)}$$

for the previous, we see that the partial derivative for each class is Gaussian. So $Q''(w) = -\sum_{i=1}^n \frac{P(\hat{y}_i = y_i | x_i, w)}{(1-P(\hat{y}_i = y_i | x_i, w))^2}$. In fact we are going to look at the discriminative model. We will assume that the prior probabilities are equal, $\pi_0 = \pi_1 = 0.5$. Then $P(\hat{y}_i = y_i | x_i, w) = \pi_i = 0.5$ when $y_i = 1$ and $P(\hat{y}_i = y_i | x_i, w) = 0.5$ when $y_i = 0$. So $Q''(w) = -\sum_{i=1}^n \frac{0.5}{(0.5)^2} = -n$. This is negative definite, so we can use gradient descent to find the minimum of the likelihood function. Can we find a closed form maximum likelihood estimation of w ?

Can't get a closed form $X^T(S^{-1})w = 0$

QDA & LDA

QDA & LDA are generative models for classification, that assume that the various classes of our data were generated by gaussian RVs. Given that we don't actually know the means or covariances of these X's, we estimate them from the data.

$$\hat{\mu}_k = \frac{1}{n_k} \sum_{i=1}^{n_k} x_i = \frac{1}{n_k} \sum_{i=1}^{n_k} (x - \hat{\mu}_k)(x - \hat{\mu}_k)^T$$

To predict which class a point belongs to, we just pick the maximum probability

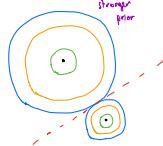
$$\begin{aligned} \hat{y} = \arg\max_k P(x|K) &= \arg\max_k \frac{P(x|K)P(K)}{K} \\ &= \arg\max_k \ln P(x|K) + \ln P(K) \\ &= \arg\max_k \ln P(x|K) \left(\hat{\mu}_K, \frac{1}{n_k} \sum_{i=1}^{n_k} x_i \right) + \ln P(K) \\ &\quad \text{Q}_K(x) \end{aligned}$$

LDA vs. QDA

LDA
• has linear boundaries

QDA
• has quadratic boundaries

Priors change the size of the level sets



Isotropic \rightarrow spherical covariance

Anisotropic \rightarrow elliptical covariance \rightarrow can be reduced to isotropic w. scaling of vars.

In LDA this leads to resulting linear boundary

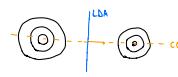
LDA, same prior:
 $\text{boundary is when equal dist from both } Q_A(x) = Q_B(x)$
 $\Rightarrow \|x - \mu_A\|^2 = \|x - \mu_B\|^2$

$$x^T (\hat{\Sigma}^{-1}(\hat{\mu}_A - \hat{\mu}_B)) + \left(\ln \frac{P(A)}{P(B)} \right) - \frac{\hat{\mu}_A^T \hat{\mu}_A - \hat{\mu}_B^T \hat{\mu}_B}{2} = 0$$

boundary still a linear function

From HW 3:

LDA and CCA are very similar: CCA on Y instead gives same result!
 \rightarrow The predicted direction of max correlation (proj. power) for CCA is the same as the normal vector to the decision boundary!



QDA, different priors:

we end up getting

$$Q_A(x) = Q_B(x)$$

\downarrow

$$\ln P(A) - \frac{1}{2}(x - \hat{\mu}_A)^T \hat{\Sigma}_A^{-1}(x - \hat{\mu}_A) = \ln P(B) - \frac{1}{2}(x - \hat{\mu}_B)^T \hat{\Sigma}_B^{-1}(x - \hat{\mu}_B)$$

unlike prior cases cannot cancel out the quadratic terms, so we are left with a quadratic boundary

\rightarrow LDA vs. CCA

Labels don't have a hard assignment anymore \rightarrow we can just look at $Q_K(x)$ as a vector of probabilities for point x for each class!

problems with simple K-means:

- no probabilistic interpretation

- each feature equally important \rightarrow related to spherical gaussians

problems with soft K-means:

- how to choose β ?

- clusters still spherical

problems with soft K-means:

- closer distances \rightarrow very skewed prob.

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

- closer points \rightarrow all other

- distance is very good

$\beta \nearrow$ more regularized with β

$\beta \searrow$ probability is $\alpha \rightarrow 0$

problems with soft K-means:

Nearest Neighbor Classification: A very intuitively simple approach to classification is just to classify each point according to its neighbours.

Classifies a locally defined (near one has w/ our own measure)



K is \Rightarrow hyperparameter that changes the decision boundary a lot

Training accuracy is always 100% (with K=1)

As $K \rightarrow \infty$ bias \Rightarrow Variance \downarrow (few K neighbors)

Bias-Variance Analysis

Let $D = \{(x_i, y_i)\}_{i=1}^n$ be our dataset and our hypothesis $h(x) = \frac{1}{K} \sum_{k=1}^K N(x_k, x, K)$ where $N(x_k, x) = \begin{cases} 1 & \text{if } x_k \text{ is one of } K \text{ closest to } x \\ 0 & \text{else} \end{cases}$

Also suppose true labels are $y_i = f(x_i) + \epsilon_i$ (basically we are applying KNN to a regression problem)

While, let x_1, \dots, x_K be the closest K , $\epsilon_i = f(x_i) + \epsilon_i$

$$\therefore \text{bias}^2 = (\mathbb{E}[h(x)] - f(x))^2 = \left(\frac{1}{K} \sum_{k=1}^K \mathbb{E}[f(x_k)] - f(x) \right)^2 = \left(\frac{1}{K} \sum_{k=1}^K f(x_k) + \mathbb{E}[\epsilon_k] \right)^2 \quad \text{when } K \rightarrow \infty \text{ bias is zero}$$

$$\therefore \text{bias}^2 = (\mathbb{E}[h(x)] - f(x))^2 = \frac{1}{K} \text{Var}\left(\frac{1}{K} \sum_{k=1}^K f(x_k) + \epsilon\right) = \frac{1}{K} \text{Var}\left(\frac{1}{K} \sum_{k=1}^K f(x_k)\right) = \frac{1}{K} \sum_{k=1}^K \text{Var}(f(x_k)) = \frac{1}{K} \text{variance}$$

Also variance = $\text{Var}(h(x)) = \frac{1}{K^2} \text{Var}\left(\frac{1}{K} \sum_{k=1}^K f(x_k) + \epsilon\right) = \frac{1}{K^2} \text{Var}\left(\sum_{k=1}^K f(x_k)\right) + \frac{1}{K^2} \sum_{k=1}^K \text{Var}(\epsilon_k) = \frac{1}{K} \sum_{k=1}^K \text{Var}(\epsilon_k)$

In contrast, $\text{Var}(h(x)) = \frac{1}{K^2} \text{Var}\left(\frac{1}{K} \sum_{k=1}^K f(x_k) + \epsilon\right) = \frac{1}{K^2} \text{Var}\left(\sum_{k=1}^K f(x_k)\right) + \frac{1}{K^2} \sum_{k=1}^K \text{Var}(\epsilon_k) = \frac{1}{K} \sum_{k=1}^K \text{Var}(\epsilon_k)$

In contrast, Linear Reg, LDA, QDA, are all parametric

Properties:

• Computational Complexity: $O(n)$ training, $O(n)$ space for storing data. Takes $O(n)$ for predictions, but active research in approximation that take less time.

• Flexibility: Can be modified to output probabilities $P(Y|X)$ by defining them as proportion of labels among the neighbors (for $K=2$)

• Can learn very complicated non-linear boundaries

• Can be applied both to discrete labels (discrete/real) and continuous (average label)

• Non-Parametric: KNN is non-parametric (does not have a fixed form) \rightarrow # parameters grows with n

(In contrast, Linear Reg, LDA, QDA, are all parametric)

• Behaviour in High dimensions: Does not do well in high dimensions, so calculating similarity usually involves

iterating all d entries of each datapoint (inner products for example take O(d) time)

Moreover, as d increases, the datapoints will tend to be further and further apart, detracting from the success of the method.

• Theoretical properties: 1-NN has impressive theoretical guarantees \rightarrow when $n \rightarrow \infty$ the expected error for 1-NN is upper bounded by $2\varepsilon^2$ where ε^2 is

The Bayes optimal error

$$\text{If } \frac{n}{K} \rightarrow \infty \text{ and } \frac{K}{n} \rightarrow 0 \text{ then Error (1NN)} \rightarrow \varepsilon^2$$

Curse of Dimensionality

Volume of a ball: as $d \rightarrow \infty$ nearly all the volume of a ball is on an arbitrarily thin shell

§ Concentration of Measure

Improving KNN

1. More data

2. Reduce dimensionality of the features / pick better features (lower some places very much)

3. Modify the distance function: Min-Max-Ki distances

$$D_p(x, z) = \left(\sum_{i=1}^d |x_i - z_i|^p \right)^{\frac{1}{p}}$$

4. Use kernels to compute distances in higher dim space with $O(d)$ computation

Sparcity: Many times we want sparse solutions. This is also a form of regularization.

Why sparsity? $\left\{ \begin{array}{l} \text{regularization} \\ \text{Speed up at test time (decreasing features)} \end{array} \right.$

Problem: $\|w\|_0$ is not convex. Solving is NP-Hard \rightarrow we will approximate this problem (in a manner a proper norm)

• LASSO: is a relaxed version of Sparse Least Squares

$$\min_w \|Xw - y\|_2^2 \quad \text{s.t. } \|w\|_1 \leq K$$

Due to strong duality, this is equivalent to

$$\min_w \|Xw - y\|_2^2 + \lambda \|w\|_1$$

Finding Optimal Solution Why CB instead of SGD?

There is no closed form solution with linear algebra due to duality (not differentiable)

However, objective is still convex, so we can use a gradient descent approach. Coordinate descent?

In the case of LASSO, our obj. fn. is strictly convex in all components of w , so we can apply coordinate descent with the guarantee of finding a global minimum.

By using coordinate descent we obtain the update:

$$w_i \leftarrow \arg\min_w L(w)$$

This update has a closed form solution. First lets rewrite the loss

$$L(w) = \|Xw - y\|_2^2 + \lambda \|w\|_1 = \left\| \frac{1}{2} (Xw - y)^T (Xw - y) + \lambda \sum_i |w_i| \right\|$$

$$= \|w_i - r\|_2^2 + \lambda |w_i| + C \quad \text{where } r = \frac{1}{2} (Xw - y)$$

Suppose $w_i > 0$: $\nabla_{w_i} L(w) = \sum_{j=1}^d 2x_{ji}(w_i x_{ji} + r_j) \Rightarrow w_i^* = \frac{-\sqrt{\frac{2}{\lambda}} r_i - \sqrt{\frac{2}{\lambda}} x_{ji}^2}{\frac{2}{\lambda} + x_{ji}^2}$

$$= \frac{-\lambda r_i - \lambda x_{ji}^2}{\lambda + 2x_{ji}^2}$$

This only holds if $w_i = \frac{-\lambda r_i}{\lambda + 2x_{ji}^2}$ is strictly positive

If negative or 0, there is no optimum in $(+, +)$

When $w_i < 0$:

$$w_i^* = \frac{\lambda x_{ji}^2}{\lambda + 2x_{ji}^2} \quad (\text{only when } \frac{\lambda x_{ji}^2}{\lambda + 2x_{ji}^2} \text{ is strictly negative})$$

If positive or 0, there is no optimum in $(-, -)$

If neither, then this is an optimum in $(-, +) \cup (+, -)$.

However we know that the most exact fit is optimum $\Rightarrow w^* = 0$ (so strictly convex)

$$\therefore w_i^* = \begin{cases} 0 & \text{if } |x_{ji}| \leq \lambda \\ \frac{\lambda x_{ji}^2}{\lambda + 2x_{ji}^2} & \text{if } |x_{ji}| > \lambda \\ \frac{-\lambda r_i - \lambda x_{ji}^2}{\lambda + 2x_{ji}^2} & \text{if } x_{ji} < 0 \end{cases} \quad \text{By increasing } \lambda \text{ make sparse more common.}$$

Not a grad. descent update? (After we fix it)

→ After having been set to 0, weights can be "restored" in subsequent iterations. [Need to center data to zero in order to parallel set features similarly]

Orthogonal Matching Pursuit

Extension of Matching Pursuit: at each iteration i, we maintain a set I^i of all features selected by the algorithm so far.

Instead of updating one component of the weight vector at a time, we update all weights corresponding to the features in I^i using least squares

Algorithm

Initialization: $w^* = 0$ and $r^* = y - Xw^*$

Initialize features: $I^0 = \emptyset$

while $\|w\|_0 < K$:

 find the feature for which length of projected residual onto x_i is minimized

$$i^* = \arg\min_i (\|r^* - \langle r^*, x_i \rangle x_i\|) = \arg\min_i \frac{|\langle r^*, x_i \rangle|}{\|x_i\|}$$

 add the feature to the set: $I^i = I^{i-1} \cup \{i\}$

 Estimate the best linear fit possible using the current features:

$$w^* = \arg\min_w \|Xw - y\|_2^2 \quad \text{where } w \in \mathbb{R}^{|I^i|}$$

 update residual vector: $r^* = Xw^* - y$

Idea: If we don't refit, we are not guaranteed that new residual will be orthogonal to the span of previous basis vectors, and therefore is not optimal.

∴ OMP ensures optimality on the set of features I^i at each step.

Consider (1) $\min_w \|r^* - \langle r^*, x_i \rangle x_i\|$ is a projection problem finding distance of the span of x_i from r^*

$$r^* = \langle r^*, x_i \rangle x_i \Rightarrow \|r^* - \langle r^*, x_i \rangle x_i\| \geq \|r^* - \langle r^*, x_i \rangle x_i\|_{\text{span}}$$

The inner optimization minimizes the distance.

$$\min_w \|\langle r^*, x_i \rangle x_i - w\|_2^2 \quad \text{where something weird here}$$

$$w = w^* + \frac{\langle r^*, x_i \rangle x_i}{\|x_i\|^2} \quad \text{why does this work?}$$

update w : $w^* = w^* + \frac{\langle r^*, x_i \rangle x_i}{\|x_i\|^2}$

update r^* : $r^* = r^* - \langle r^*, x_i \rangle x_i$

Why does this work?

Consider (1) $\min_w \|r^* - \langle r^*, x_i \rangle x_i\|$ is a projection problem finding distance of the span of x_i from r^*

$$r^* = \langle r^*, x_i \rangle x_i \Rightarrow \|r^* - \langle r^*, x_i \rangle x_i\| \geq \|r^* - \langle r^*, x_i \rangle x_i\|_{\text{span}}$$

The inner optimization minimizes the distance.

$$\min_w \|\langle r^*, x_i \rangle x_i - w\|_2^2 \quad \text{where something weird here}$$

$$w = w^* + \frac{\langle r^*, x_i \rangle x_i}{\|x_i\|^2}$$

$$= w^* + \frac{\langle r^*, x_i \rangle x_i}{\|x_i\|^2} - \frac{\langle r^*, x_i \rangle x_i}{\|x_i\|^2} + \frac{\langle r^*, x_i \rangle x_i}{\|x_i\|^2}$$

$$= w^* + \frac{\langle r^*, x_i \rangle x_i}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{\|x_i\|^2}$$

$$= \arg\min_w \frac{\|\langle r^*, x_i \rangle x_i - w\|_2^2}{$$

Decision Trees • Can be used for both classification & regression
• Each branch represents a split in the data
• Can represent an extremely complex classifier.

Training
Trained in recursive fashion, going downwards from the root, at each step finding the best split of the data.
The base case is when we stop splitting the data, and give 2 prediction (leaf nodes).

In CS193 we only consider splits of the type $x_1 \leq v$, $x_2 \leq v$. These are finitely many possible splits, and they can be found by ordering the x datapoints by x_1 .

Loss: At each split, we want to reduce the classifier's uncertainty, as much as possible!
how to quantify?

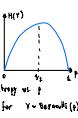
Entropy & Information

"Surprise" of a discrete R.V. Y taking on value y_i :

$$\log \frac{1}{P(Y=y_i)} = -\log P(Y=y_i)$$

Entropy of an R.V. Y is the expected surprise score:

$$H(Y) = \mathbb{E}[-\log P(Y)] = -\sum_k P(Y=k) \log P(Y=k)$$



In general, entropy is higher the closer the distribution is to a uniform distribution.

Entropy Approach to Decision Trees

Concretely we don't have the probability for each class, but we can estimate it as:

$$P(Y=k) = \frac{\# y=k}{n}$$

We then know that we want to pick a split that minimizes entropy: of what?

$$\text{point w.r.t. } \begin{cases} \text{one side of } v \\ H(Y|X_{v>0}) \end{cases} \quad \begin{cases} \text{other side} \\ H(Y|X_{v<0}) \end{cases}$$

We want natural $H(Y)$ but can't afford the other two.

Assume we want to do a minimize the "conditional entropy":

$$\min H(Y|X_{v>0}) = P(X_{v>0}=1)H(Y|X_{v>0}) + P(X_{v>0}=0)H(Y|X_{v>0}=0)$$

Equivalently, we can maximize the mutual information (information gain):

$$\max I(X_{v>0}; Y) = H(Y) - H(Y|X_{v>0})$$

↳ always non-negative

Random Forests

The main idea is to combine multiple Decision Tree classifiers. This technique falls under **ensembling**.

Ensembling relies on the idea that having multiple partially-independent measurements will reduce the variance of our prediction, without biasing the estimate.

Because the way we created decision trees was deterministic, we need to introduce stochasticity in order to create a forest of different trees.

- **Bagging (Bootstrap Aggregating):** sample m datapoints uniformly with replacement, and use them as training set
- **Feature Randomization:** at each split, only consider $K < d$ random features

Both the size of random subsample set number of features can be chosen with cross-validation.

Boosting

Random Forests: All consider same importance with subtrees by accuracy

Boosting: Consider models (weak learners) in a principled manner

Boosting Idea: Combine models (weak learners) in a principled manner

such that the overall learner correctly does not predict very well

In fact, overall predictor is an additive combination of pieces which are selected one-by-one in a greedy fashion

Consider also pruning: resulting splits if reduces validation error

Boosting vs Matching Pursuit

In both, overall predictor is an additive combination of pieces which are selected one-by-one in a greedy fashion

Notice the above are not mutually exclusive

Consider also pruning: resulting splits if reduces validation error

AdaBoost (Adaptive Boosting)

One of the most common boosting algorithms.

Assume we have a dataset $\{(x_i, y_i)\}_{i=1}^n$, $x_i \in \mathbb{R}^d$, $y_i \in \{-1, +1\}$

Algorithm

i) Initialize weights w_1, \dots, w_n for training points

ii) Repeat for $m=1, \dots, M$

 a) Build a classifier $G_m: \mathbb{R}^d \rightarrow \{-1, +1\}$ where datapoints are weighted according to w_i

 b) Compute weighted error rate $e_m = \frac{\sum_i w_i G_m(x_i)}{\sum_i w_i}$

 c) Re-weight the training points: $w_i \leftarrow w_i \begin{cases} \frac{e_m}{1-e_m} & \text{if misclassified by } G_m \\ \frac{1-e_m}{e_m} & \text{else} \end{cases}$

 d) Optional: normalize weights w_i to sum to 1

Derivation

The derivation will be intuitive. Suppose we have completed the classifiers G_1, \dots, G_{m-1} along with corresponding weights w_1, \dots, w_{m-1} . We now want to complete the next classifier G_m and its weight e_m . The output of the model is for $L(y, f_m(x)) = \frac{1}{n} \sum_i L(y_i, f_m(x_i))$. We want to minimize empirical risk:

$$\min_{G_m} \frac{1}{n} \sum_i L(y_i, f_m(x_i) + \epsilon_m G_m(x_i))$$

For AdaBoost, we use exponential loss: $L(y, f) = e^{-yf}$

$$\Rightarrow \min_{G_m, \epsilon_m} \frac{1}{n} \sum_i e^{-y_i(f_m(x_i) + \epsilon_m G_m(x_i))}$$

$$\Rightarrow \min_{G_m, \epsilon_m} \frac{1}{n} \sum_i e^{-y_i f_m(x_i)} e^{-\epsilon_m G_m(x_i)}$$

$$= \arg \min_{G_m, \epsilon_m} \sum_i w_i^{(m)} e^{-y_i f_m(x_i)} e^{-\epsilon_m G_m(x_i)}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} \sum_i w_i^{(m)} e^{-y_i f_m(x_i)} e^{-\epsilon_m G_m(x_i)}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$

$$= \arg \min_{\epsilon_m, w_i^{(m)}} (1 - \epsilon_m)^{w_i^{(m)}} e^{-\epsilon_m} e^{w_i^{(m)}}$$