

PART 1

Linear Regression

x_0	x_1	y
1	1	2
1	4	3
1	6	4

$$y = w_0 + w_1 x_1$$

$$= \bar{w}^T \bar{x}$$

Mean Squared
error

y observed	\hat{y} predicted
2	$-(w_0 + w_1)^2$
3	$-(w_0 + 4w_1)^2$
4	$-(w_0 + 6w_1)^2$

Dummy Variable

$$\bar{x} = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} \quad \bar{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$

$$\nabla J = \left[\frac{\partial J}{\partial w_0}, \frac{\partial J}{\partial w_1} \right]^T = \bar{0} \rightarrow \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\Rightarrow w_0 = 1, w_1 = 0.5$$

Matrix Notation

$$\begin{array}{c} \overline{x}_0 \quad \overline{x}_1 \quad \overline{y} \\ \hline \boxed{\begin{matrix} 1 & 1 \\ 1 & 4 \\ 1 & 6 \end{matrix}} \quad \boxed{\begin{matrix} 2 \\ 3 \\ 4 \end{matrix}} \end{array} \quad \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} \quad 2 \times 1$$

$\bar{X} \quad 3 \times 2 \quad \bar{y} \quad 3 \times 1$

$$J(\bar{w}) = \frac{1}{3} (\bar{y} - \bar{X}\bar{w})^T (\bar{y} - \bar{X}\bar{w})$$

$$\begin{matrix} 3 \times 1 & 3 \times 2 & 2 \times 2 \\ & \underbrace{\quad\quad\quad}_{3 \times 1} & \\ \downarrow & \downarrow & \downarrow \end{matrix}$$

Design Matrix

$$\nabla J = 0$$

$$\Rightarrow \bar{w} = \underbrace{(\bar{X}^T \bar{X})^{-1}}_{2 \times 2} \bar{X}^T \bar{y} \quad \underbrace{\begin{matrix} 2 \times 3 \\ 3 \times 1 \end{matrix}}_{2 \times 1}$$

Ordinary least square (OLS)
Solution

probabilistic Linear Regression

So far: $\bar{y} = \bar{w}^T \bar{x}$ \bar{y} completed determined by \bar{x}

More reasonable model

$$\underline{\bar{y}} = \bar{w}^T \bar{x} + \text{impact of other minor factors}$$

$$= \bar{w}^T \bar{x} + \underline{\varepsilon} \quad \varepsilon \sim N(0, \sigma^2)$$

Noise ε ~ Normal distribution

$$P(\bar{y} | \bar{x}, \theta) = N(\mu(\bar{x}), \sigma^2)$$

$$\theta = (\bar{w}, \sigma)$$

$$\bar{w}^T \bar{x}$$

$$\mathcal{D} = \{\bar{x}_i, y_i\}_{i=1}^N$$

Minimize

$$-\frac{1}{N} \sum_{i=1}^N \log p(y_i | \bar{x}_i, \Theta)$$

$$= -\frac{1}{N} \sum_{i=1}^N \log \left[\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y_i - \bar{w}^T \bar{x}_i)^2}{2\sigma^2}} \right]$$

$$= -\frac{1}{N} \sum_{i=0}^N \log \frac{1}{\sqrt{2\pi}\sigma} + \frac{1}{2\sigma^2} \frac{1}{N} \sum_{i=1}^N (y_i - \bar{w}^T \bar{x}_i)^2$$

$$\sigma = 1$$

MSE



Polynomial Regression

$$\begin{array}{c|ccc} x_0 & x_1 & y \\ \hline 1 & 1 & 2 \end{array}$$

$$\begin{array}{c|cc} & 4 & 3 \\ \hline 1 & 6 & 4 \end{array}$$

$$\begin{array}{c|cccc} & x_2 & x_3 & & \\ \hline x_0 & x_1 & x^2 & x^3 & y \\ \hline 1 & 1 & 1 & 1 & 2 \\ 4 & 16 & 64 & 3 & \\ 6 & 36 & 216 & 4 & \end{array}$$

$$y = w_0 + w_1 x_1 = \bar{w}^T \bar{x}$$

More complex model

$$y = w_0 + w_1 x_1 + w_2 x_1^2 + w_3 x_1^3$$

$$\bar{x} = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} \Rightarrow \phi(\bar{x}) = \begin{bmatrix} x_0 \\ x_1 \\ x_1^2 \\ x_1^3 \end{bmatrix}$$

↑
feature transformation
mapping

Poly Reg. = Poly feature transformation
+ Linear Reg

Q: What polynomial should be used?

order d

1

3

10

Hypothesis Space

$$Y = w_0 + w_1 X_1$$

$$S_1 = \{ Y = w_0 + w_1 X_1 \mid w_0, w_1 \in \mathbb{R} \}$$

$$\bar{w}_1 = (w_0, w_1)$$

$$Y = 1 + X_1$$

$$Y = -0.2 + 5X_1$$

$$S_2 = \{ Y = w_0 + w_1 X_1 + w_2 X_1^2 + w_3 X_1^3 \mid w_0, w_1, w_2, w_3 \in \mathbb{R} \}$$

$$S_1 \subset S_2$$

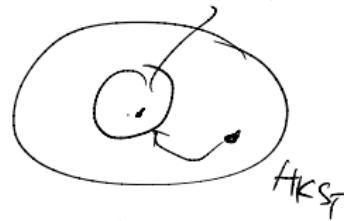
$$\bar{w}_2 = (w_0, w_1, w_2, w_3)^T$$

52/12

$$\min_{\bar{w}_1} J(\bar{w}_1)$$

\geq

$$\min_{\bar{w}_2} J(\bar{w}_2)$$



HKUST

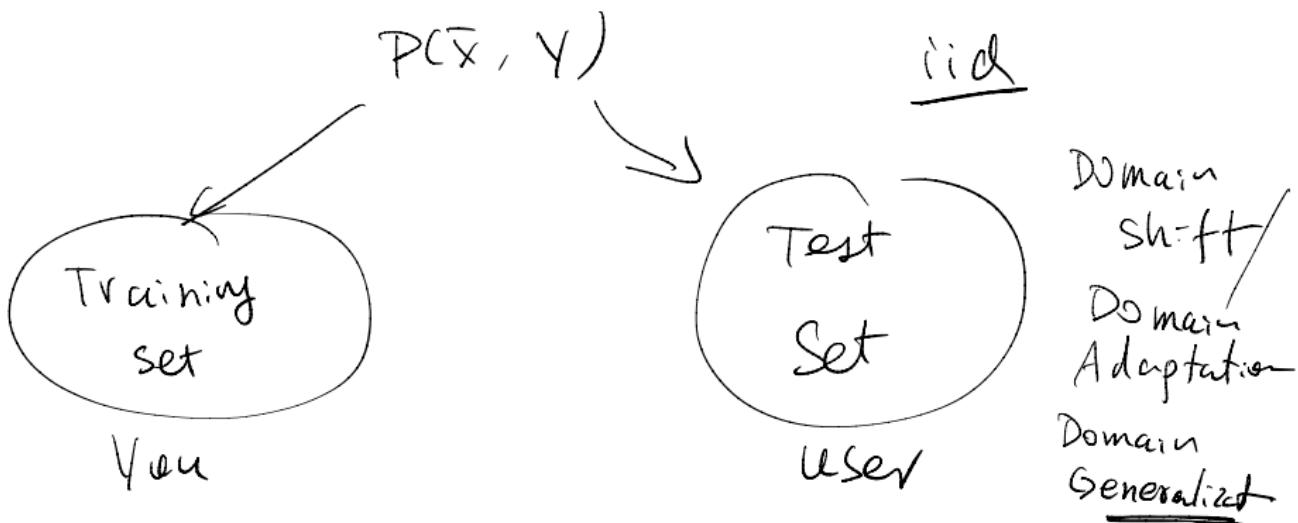
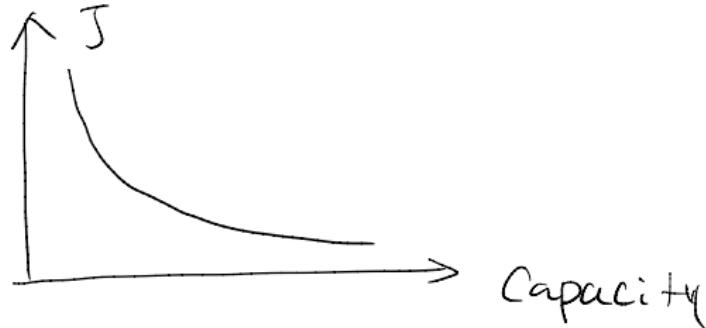
height of shortest

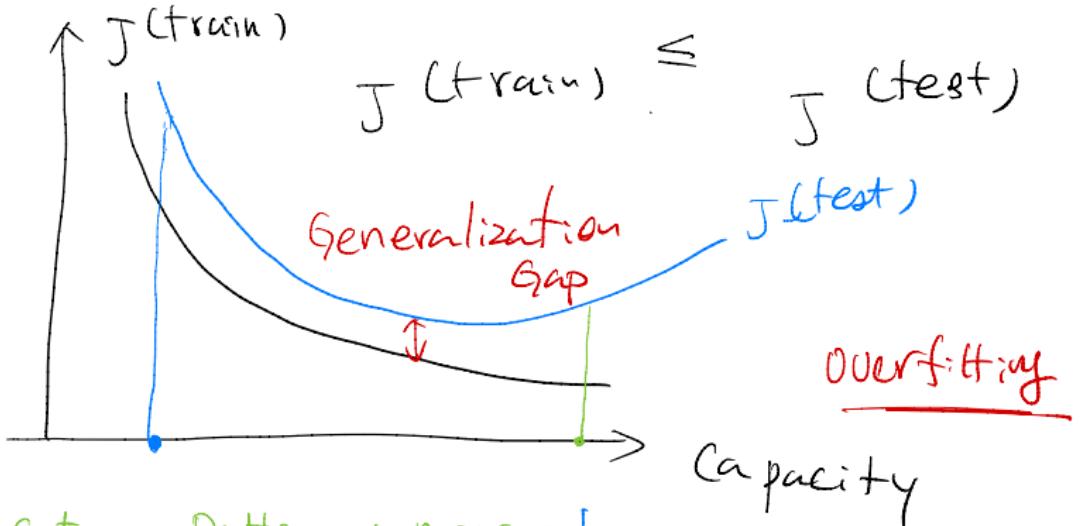
person in Comp5212

\geq (B)

height of shortest

person in HKUST





Training Set: Pattern + noise₁
 Test Set: Pattern + noise₂

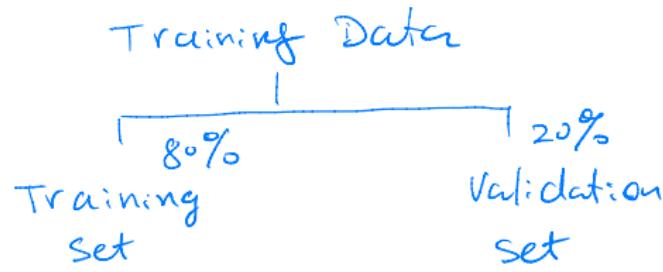
Model too simple
 * cannot fit pattern
 * $J(\text{train}), J(\text{test})$ large

Model too complex
 * Fit pattern + noise₁
 * $J(\text{train})$ small
 * $J(\text{test})$ large
 dist(Noise₁, Noise₂)

Underfitting

2022-02-18

Validation for choosing hyperparameters d



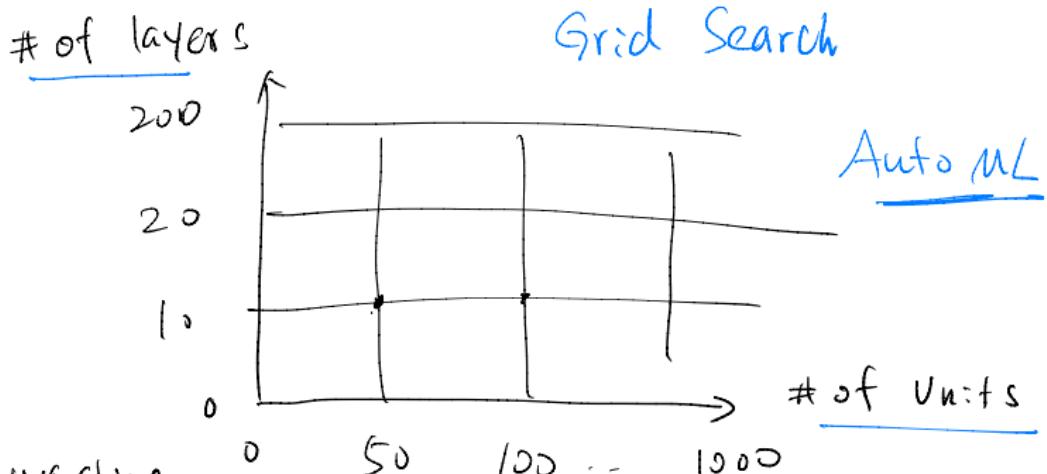
$d = 1 \quad 1.5 \quad 1.6$

$d = 5 \quad 0.3 \quad 0.35 \quad \checkmark$

$d = 10 \quad 0 \quad 0.5$

what values to consider for d ?

1, 2, 3, -- -

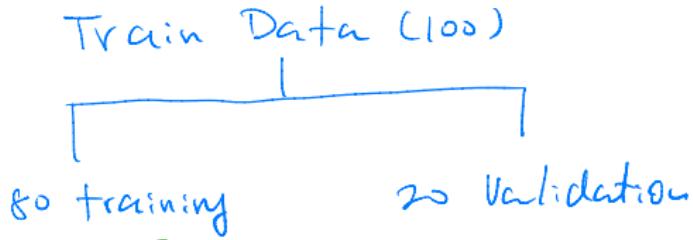


Exhaustive

200,000 models

1

when data is scarce



- * Not use all data in model construction
- * Validation error has high Variance

Income of HKer :

Method 1: Survey 5 people

Average

Method 2: Survey 1000 people

Average

Team 1	Team 2
12k	21k

$$\bar{x}_1 = \frac{1}{5} - 0 /$$

$$\bar{x}_2 = \frac{18}{1000} - 0 /$$

Regularization

$$\min J(\bar{w})$$

$$= \frac{1}{N} \sum_{i=1}^N (y_i - (w_0 + \bar{w}^\top \phi(\bar{x}_i)))^2$$

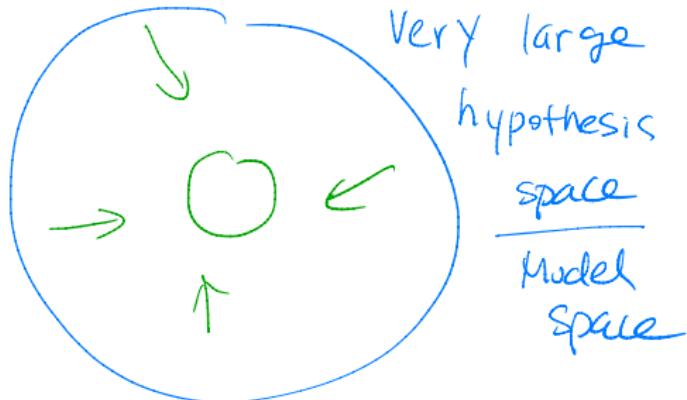
$$+ \lambda \|\bar{w}\|_2^2$$

Two competing objective

* min loss

* min weights

Bias
↓
 $w = \begin{bmatrix} w_1 \\ \vdots \\ w_D \end{bmatrix}$



Prior: Simpler model

$$w_1^2 + w_2^2 + \dots + w_D^2$$

pushing down the weights
for the feature
unimportant

$$Y = w_0 + w_1 x_1 + w_2 x_1^2 + \cancel{w_3 x_1^3} + \cancel{w_4 x_1^4}$$

$$0.8 \quad 1.5 \quad 0.2 \quad 0.001 \quad 0.0001$$

$$J(\bar{w}) = \frac{1}{N} \sum_{i=1}^N (y_i - (w_0 + \bar{w}^\top \phi(x_i)))^2 + \lambda \|\bar{w}\|_2^2$$

λ : balance parameter

$$\boxed{\cancel{+} \cancel{\lambda^2}}$$

$$\lambda = 0$$

How to choose λ ?

$$\lambda = 10$$

$$\lambda = 1000$$

$$J(\bar{w}) = \frac{1}{N} \sum_{i=1}^N (y_i - (w_0 + \bar{w}^\top \phi(\bar{x}_i))^2 + \lambda \|\bar{w}\|_2^2$$

$w_1^2 + w_2^2 + \dots$
 \nwarrow
 $\lambda \|\bar{w}\|_2^2$

L₂ regularization, weight decay, ridge regression

$$J(\bar{w}) = \frac{1}{N} \sum_{i=1}^N (y_i - (w_0 + \bar{w}^\top \phi(\bar{x}_i))^2 + \lambda \|\bar{w}\|_1$$

$|w_1| + |w_2| + \dots$
 \nwarrow

L₁ regularization

LASSO (least absolute shrinkage
and selection operation)

Sparse Model

$$J(\bar{w}) = l(\bar{w}) + r(\bar{w}) \quad \text{minimized at } w^*$$

$$a = l(\bar{w}^*), b = r(\bar{w}^*)$$

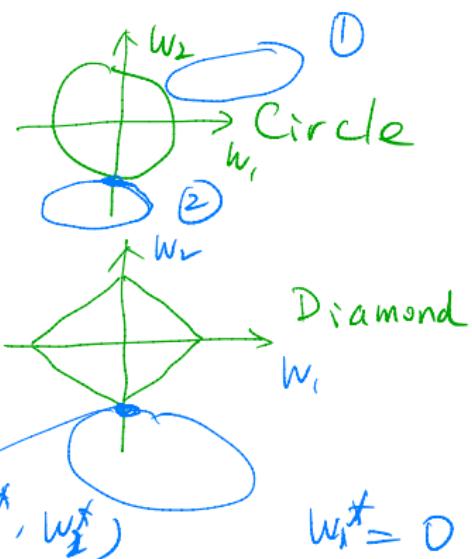
$$l(\bar{w}) = a \quad \text{Contour ①}$$



$$r(\bar{w}) = b \quad \text{Contour ②}$$

$$\begin{cases} L_2: \quad w_1^2 + w_2^2 = b \\ L_1: \quad |w_1| + |w_2| = b \end{cases}$$

w^* is where ① & ② meet



Regularization

$$\min \ell(D, Q) + \lambda f(Q)$$

Encourage Q to have certain properties

$$\ell(D, Q) + \lambda \|W\|_2^2 \quad \text{Simple}$$

$$\ell(D, Q) + \lambda I(z, z_r) \quad \text{Disentanglement}$$

$\nwarrow \nwarrow$ latent variables in model

$$\ell(D, Q) - \lambda H[Q(Y|X)] \quad \text{reduce confidence}$$

0.99
Max entropy

of result

classes ! . . . -

Linear Regression: $\{\bar{x}_i, y_i\}_{i=1}^N \quad y_i \in \mathbb{R}^1$

$$\bar{x} \longrightarrow \bar{w}^T \bar{x} \quad P(y|\bar{x}) = N(\bar{w}^T \bar{x}, \sigma)$$

Logistic Regression: $\{\bar{x}_i, y_i\}_{i=1}^N \quad y_i \in \{0, 1\}$

Want : $P(y|\bar{x}, \bar{w}) \in [0, 1]$

$\bar{x} \rightarrow \bar{w}^T \bar{x} \in \mathbb{R}^1$

$$P(y=1|\bar{x}, \bar{w}) = \sigma(\bar{w}^T \bar{x}) = \frac{1}{1 + e^{-\bar{w}^T \bar{x}}}$$

$$P(y=0|\bar{x}, \bar{w}) = 1 - \frac{1}{1 + e^{-\bar{w}^T \bar{x}}} = \frac{e^{-\bar{w}^T \bar{x}}}{1 + e^{-\bar{w}^T \bar{x}}}$$

$$P(Y=1|\bar{x}, \bar{w}) = \sigma(\bar{w}^T \bar{x}) = \frac{1}{1 + e^{-\bar{w}^T \bar{x}}}$$

$$P \in [0, 1] \quad P(Y=0|\bar{x}, \bar{w}) = 1 - \frac{1}{1 + e^{-\bar{w}^T \bar{x}}} = \frac{e^{-\bar{w}^T \bar{x}}}{1 + e^{-\bar{w}^T \bar{x}}}$$

$$\text{logit}(P) = \log \frac{P}{1-P}$$

$\hookrightarrow \in (-\infty, \infty)$

$$\begin{aligned}\text{Logit}[P(Y=1|\bar{w}, \bar{x})] &= \log \frac{P(Y=1|\bar{x}, \bar{w})}{P(Y=0|\bar{x}, \bar{w})} \\ &= \log e^{\bar{w}^T \bar{x}} \\ &= \bar{w}^T \bar{x}\end{aligned}$$

logistic regression : linear regression for logit

$$\hat{Y} = 1 \iff P(Y=1 | \bar{x}, \bar{w}) > P(Y=0 | \bar{x}, \bar{w})$$

$$\iff \frac{P(Y=1 | \bar{x}, \bar{w})}{P(Y=0 | \bar{x}, \bar{w})} > 1$$

$$\iff e^{\bar{w}^T \bar{x}} > 1$$

$$\bar{X} = [x_1, x_2, \dots, x_D]$$

$$\iff \bar{w}^T \bar{x} > 0 \quad \underbrace{w_0 + w_1 x_1 + \dots + w_D x_D}_{> 0}$$

Logistic Regression is a

Linear classifier

To learn \bar{w} , minimize cross entropy $l(\bar{w})$

$$\begin{aligned} J(\bar{w}) &= -\frac{1}{N} \sum_{i=1}^N \log P(Y_i | \bar{x}_i, \bar{w}) \\ &= -\frac{1}{N} \sum_{i=1}^N [Y_i \log P(Y_i=1|...) + (1-Y_i) \log (1-P(Y_i=1|...))] \\ &= -\frac{1}{N} \sum_{i=1}^N [Y_i \log \sigma(\bar{w}^T \bar{x}_i) + (1-Y_i) \log (1-\sigma(\bar{w}^T \bar{x}_i))] \end{aligned}$$

\Rightarrow Binary Cross Entropy

$$\frac{\log P(Y=1|...)}{\log P(Y=0|...)} = \frac{Y \log P(Y=1|...) + (1-Y) \log (1-P(Y=1|...))}{0 + \log (1-P(Y=1|...))} = \log P(Y=1|...)$$

$$Y=0 \quad \log P(Y=0|...) \quad 0 + \log (1-P(Y=1|...)) = \log P(Y=0|...)$$

$$Y=1 \quad \log P(Y=1|...) \quad \log P(Y=1|...) + 0$$

Gradient Descent

$$J'(w) = \frac{d J(w)}{d w} = \lim_{\varepsilon \rightarrow 0} \frac{J(w + \varepsilon) - J(w)}{\varepsilon}$$

For small ε , $J'(w) \approx \frac{J(w + \varepsilon) - J(w)}{\varepsilon}$

Approximate J around w

$$J(w + \varepsilon) = J(w) + \varepsilon J'(w) \quad (1)$$

Reduce J , how to change w according to (1)?

* $J'(w) > 0$, $\varepsilon < 0$, decrease w

* $J'(w) < 0$, $\varepsilon > 0$, increase w

$$w \leftarrow w - \alpha J'(w)$$

\nwarrow step size / learning rate

$$\bar{w} = (w_0, w_1, \dots, w_D)^T \quad \bar{x} = (1, x_1, x_2, \dots, x_D)^T$$

$$J(\bar{w}) = -\frac{1}{N} \sum_{i=1}^N [y_i \underbrace{\log \sigma(\bar{w}^T \bar{x}_i)} + (1-y_i) \underbrace{\log(1-\sigma(\bar{w}^T \bar{x}_i))}]$$

i-th example $\bar{x}_i = (x_{i,0}, x_{i,1}, \dots, x_{i,D})$

$$z = \bar{w}^T \bar{x} = w_0 x_0 + w_1 x_1 + \dots + w_D x_D$$

$$\frac{\partial z}{\partial w_1} = x_1$$

$$\frac{\partial \bar{w}^T \bar{x}}{\partial w_j} = x_j$$

$$\frac{d\sigma(z)}{dz} = - \frac{1}{(1+e^z)^2} \quad \frac{d(1+e^{-z})}{dz} \quad \sigma(z) = \frac{1}{1+e^{-z}}$$

$$= - \frac{1}{(1+e^{-z})^2} \quad (- \frac{e^{-z}}{(1+e^{-z})})$$

$$= \frac{1}{(1+e^{-z})} \quad \frac{\underline{e^{-z} + 1} - 1}{(1+e^{-z})}$$

$$= \sigma(z) \left(1 - \frac{1}{1+e^{-z}} \right)$$

$$= \sigma(z) (1 - \sigma(z))$$

$$\frac{\partial}{\partial w_j} [y \log \sigma(z) + (1-y) \log (1-\sigma(z))]$$

$$= \frac{\partial}{\partial z} [y \log \sigma(z) + (1-y) \log (1-\sigma(z))] \frac{\partial z}{\partial w_j}$$

$$= \left[y \frac{1}{\sigma(z)} \frac{d \sigma(z)}{dz} + (1-y) \frac{1}{1-\sigma(z)} \frac{d(1-\sigma(z))}{dz} \right] x_j$$

$$= \left[y \cancel{\frac{1}{\sigma(z)}} \sigma(z) (1-\sigma(z)) - (1-y) \cancel{\frac{1}{1-\sigma(z)}} \sigma(z) (1-\sigma(z)) \right] x_j$$

$$= [y (1-\sigma(z)) - (1-y) \sigma(z)] x_j$$

$$= [y - y\sigma(z) - \sigma(z) + y\sigma(z)] x_j$$

$$= [y - \sigma(z)] x_j$$

Update Rule

$$w_j \leftarrow w_j + \alpha \frac{1}{N} \sum_{i=1}^N [y_i - \sigma(\mathbf{w}^\top \mathbf{x}_i)] x_{i,j}$$

observed predicted

$x_{i,j} > 0$, $y_i > \hat{y}_i$ increase \hat{y}_i , increase w_j

$$[y_i - \sigma(\cdot)] x_{ij} > 0$$

$$\forall \alpha \in \Gamma \exists \beta \in \Gamma \exists \gamma \in \Gamma \exists \delta \in \Gamma \exists \epsilon \in \Gamma$$

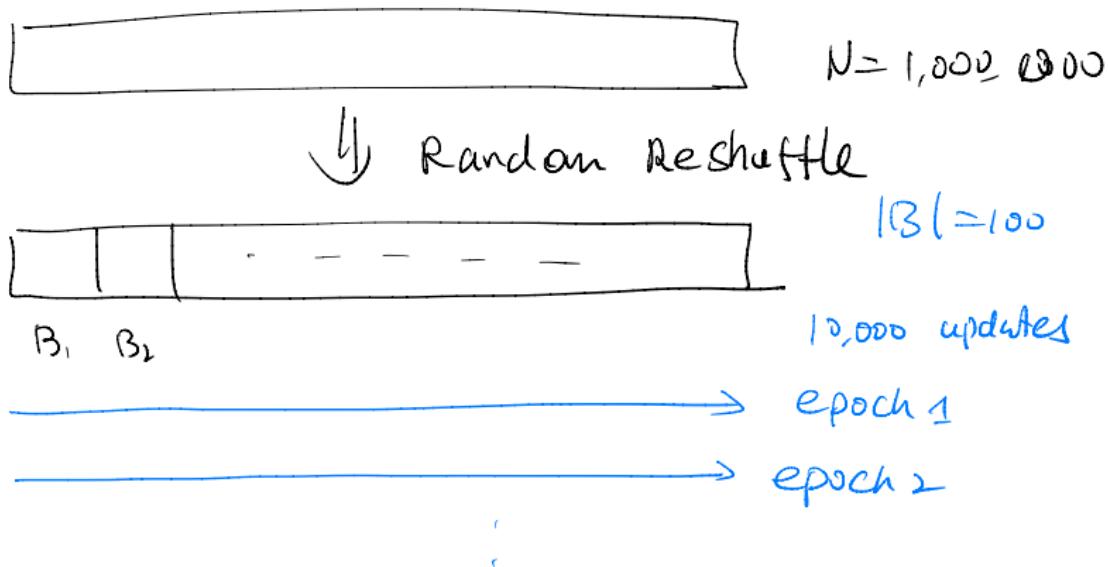
Decrease γ_i , decrease w_j

$$[Y_{i-\sigma(-j)} x_{ij} < 0]$$

$x_{i,j} < 0$, $y_i > \hat{y}_i$, increase γ , decrease w_j

$$[Y_i - \sigma(\cdot)] x_{ij} < 0$$

Stochastic Gradient Descent



Softmax Regression

$$\{\bar{x}_i, y_i\}_{i=1}^N$$

$$y_i = \{1, 2, \dots, c\}$$

↑

10

\bar{x} which class?

$$\bar{w}_1 \cdot \bar{x} e^{\bar{w}_1^T \bar{x}} / Z = P(Y=1| \cdot) \in [0, 1]$$

$$\bar{x} \quad \bar{w}_2 \cdot \bar{x} \quad e^{\bar{w}_2^T \bar{x}} / Z = P(Y=2| \cdot) \in [0, 1]$$

:

:

$$\bar{w}_c \cdot \bar{x} e^{\bar{w}_c^T \bar{x}} / Z = P(Y=c| \cdot) \in [0, 1]$$

$$Z = \sum_{c=1}^C e^{\bar{w}_c^T \bar{x}}$$

$C=2$, equivalent
to logistic Regression

$$\bar{w}_c = (w_{c,0}, \dots, w_{c,D})$$

$$\bar{W} = (\bar{w}_1, \dots, \bar{w}_C)$$

2022-02-25

Newton's Method

Task : find w s.t. $J(w) = 0$

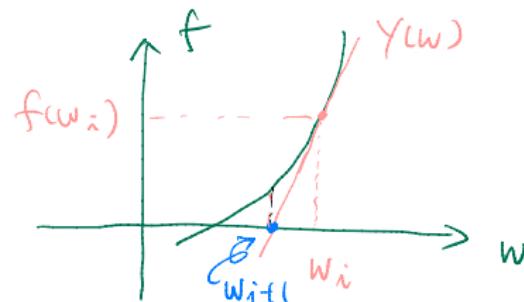
Let $f(w) = J(w)$

Find w s.t. $f(w) = 0$

* Pick w_0

* Iterate :

$$w_{i+1} \leftarrow w_i - \frac{f(w_i)}{f'(w_i)}$$



Approximate f using linear func

$$Y(w) = f'(w_i)w + f(w_i) - f'(w_i)w_i$$

$$\begin{aligned} Y(w_i) &= f'(w_i)w_i + b \\ &= f(w_i) \\ &\stackrel{=} 0 \end{aligned}$$

$$w_{i+1} \leftarrow w_i - \frac{J'(w_i)}{J''(w_i)} \leftarrow \text{2nd order}$$

GD

$$w_{i+1} \leftarrow w - \alpha J'(w_i)$$

$$b = f(w_i) - f'(w_i)w_i$$

$$\bar{w} = [w_0, w_1, \dots, w_D]^T$$

$$\nabla J = \left[\frac{\partial J}{\partial w_0}, \frac{\partial J}{\partial w_1}, \dots, \frac{\partial J}{\partial w_D} \right]^T \quad \text{Gradient vector}$$

Hessian Matrix

$$H = \begin{bmatrix} \frac{\partial^2 J}{\partial w_0 \partial w_0} & \cdots & \frac{\partial^2 J}{\partial w_0 \partial w_D} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 J}{\partial w_D \partial w_0} & \cdots & \frac{\partial^2 J}{\partial w_D \partial w_D} \end{bmatrix}$$

Newton's update Rule

$$\bar{w} \leftarrow \bar{w} - H^{-1} \nabla J \quad \underbrace{\qquad}_{O((D+1)^{2.7})}$$

Best matrix Inversion algo

$$[f_0(\bar{w}), \dots, f_D(\bar{w})]^T$$

Jacobian Matrix

$$\begin{bmatrix} \frac{\partial f_0}{\partial w_0} & \cdots & \frac{\partial f_0}{\partial w_D} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_D}{\partial w_0} & \cdots & \frac{\partial f_D}{\partial w_D} \end{bmatrix}$$

Hessian matrix is Jacobian matrix
of Gradient Vector

Taylor expansion

$$J(w+\varepsilon) = J(w) + J'(w)\varepsilon + \frac{1}{2!} J''(w)\varepsilon^2 + \frac{1}{3!} J'''(w)\varepsilon^3 + \dots$$

1st Order approximation

$J(w+\varepsilon) \approx J(w) + J'(w)\varepsilon \Rightarrow$ Gradient Descent
Algo

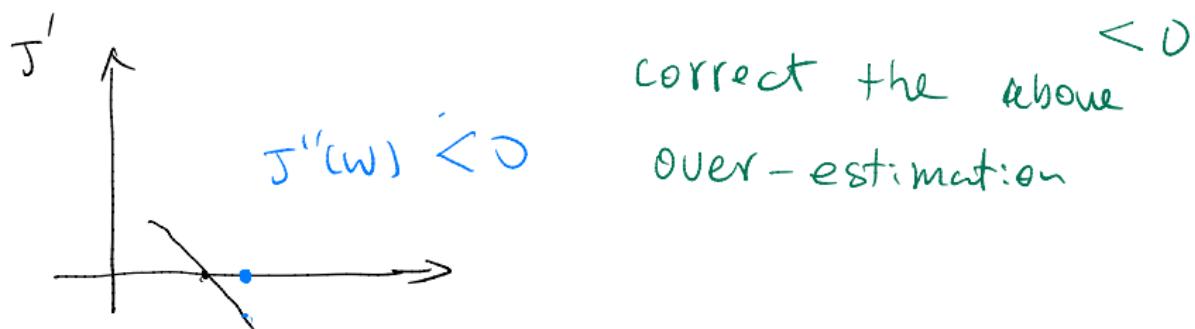
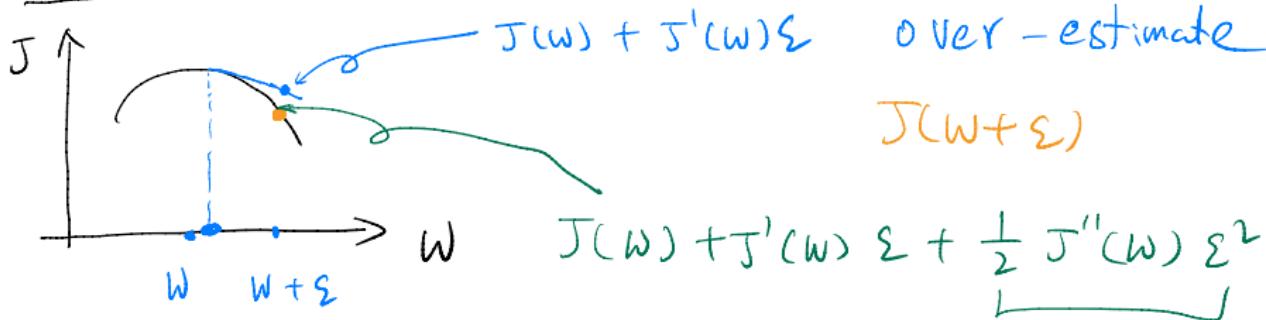
2nd order approximation

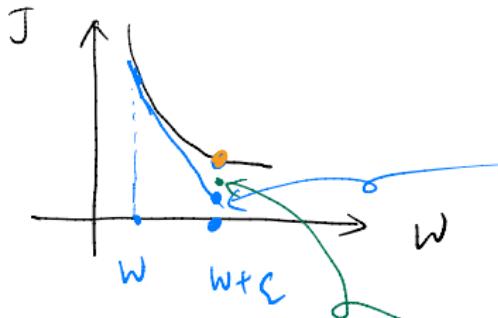
$$J(w+\varepsilon) \approx \underbrace{J(w) + J'(w)\varepsilon + \frac{1}{2} J''(w)\varepsilon^2}_{f(\varepsilon)} = f(\varepsilon)$$

$$\frac{df}{d\varepsilon} = J'(w) + J''(w)\varepsilon = 0 \quad \text{Newton's}$$

$$\Rightarrow \varepsilon = -\frac{J'(w)}{J''(w)} \quad \text{Method}$$

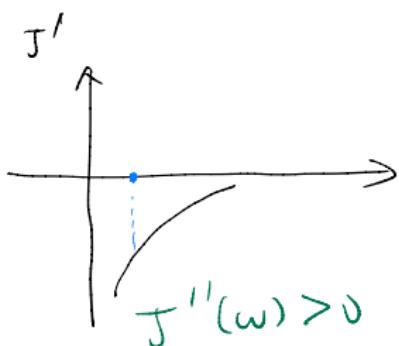
Advantage of Newton's Method





$J(w) + J'(w) \Sigma$
under-estimate

$$J(w + \zeta)$$



$$J(w) + J'(w) \Sigma + \frac{1}{2} \underline{J''(w)} \Sigma^2$$

Corrects the above
under estimation

optimization Approach

$$D = \{ \bar{x}_i, y_i \}_{i=1}^N \quad y_i \in \{-1, +1\}$$

Linear classifier : $\hat{y} = \text{sign}(\underbrace{w^T \bar{x} + b}_z) = \begin{cases} +1 & z > 0 \\ 0 & z = 0 \\ -1 & z < 0 \end{cases}$

zero/one loss

$$L_{0/1}(y_i, \hat{y}_i) = \begin{cases} 1 & \hat{y}_i \neq y_i \\ 0 & \hat{y}_i = y_i \end{cases}$$

$$L(\bar{w}, b) = \frac{1}{N} \sum_{i=1}^N L_{0/1}(y_i, \hat{y}_i) \quad \text{Empirical loss}$$

Objective $\arg \min_{\bar{w}, b} L(\bar{w}, b)$

$$\hat{Y} = \text{sign}(W^T \bar{X} + b) = \begin{cases} +1 & z > 0 \\ 0 & z = 0 \\ -1 & z < 0 \end{cases} \quad Y \in \{-1, +1\}$$

$$L_{0/1}(y_i, \hat{y}_i) = \begin{cases} 1 & \hat{y}_i \neq y_i \\ 0 & \hat{y}_i = y_i \end{cases}$$

$$\frac{1}{2} \sum_i (y_i z_i - \hat{y}_i)^2$$

observed predicted

$$= \begin{cases} 1 & \left\{ \begin{array}{l} y_i = 1, z_i \leq 0, \hat{y}_i = 0, -1 \\ y_i = -1, z_i \geq 0, \hat{y}_i = 1, 0 \end{array} \right\} y_i \neq \hat{y}_i \\ 0 & \end{cases}$$

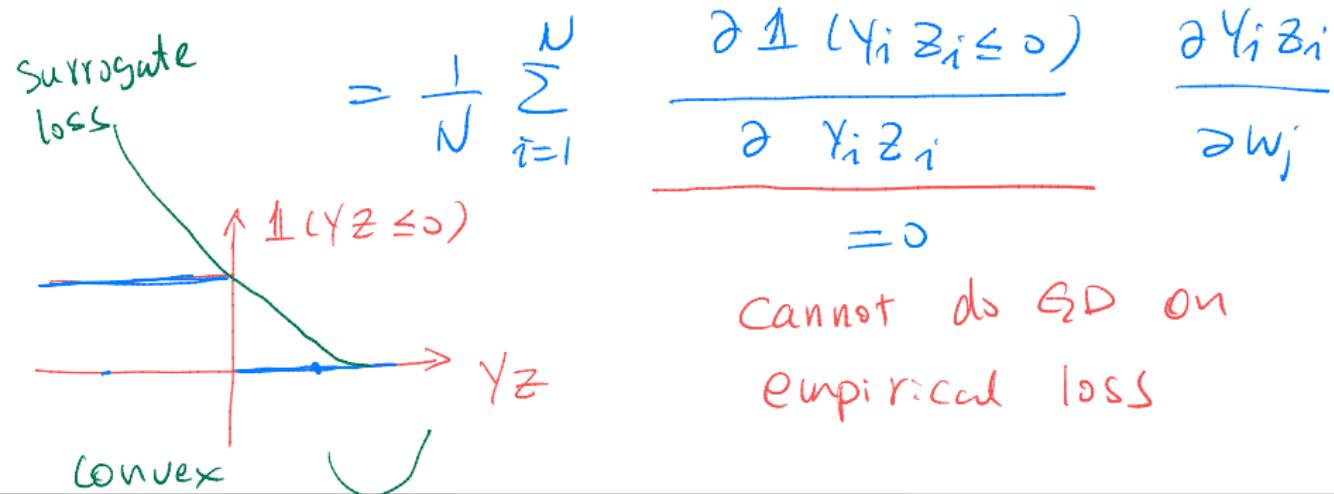
$$= \angle_{0/2} (y_i, \hat{y}_i)$$

$$Y_i = \hat{Y}_i$$

$$L(\bar{w}, b) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}(y_i z_i \leq 0)$$

$$\bar{w} = [w_1, \dots, w_d]^T$$

$$\frac{\partial L}{\partial w_j} = \frac{1}{N} \sum_{i=1}^N \frac{\partial \mathbb{1}(y_i z_i \leq 0)}{\partial w_j}$$



Want

$$\min \frac{1}{N} \sum_{i=1}^N L_{\text{ols}}(y_i, \hat{y}_i)$$

training

error

Do

$$\min \frac{1}{N} \sum_{i=1}^N L_{\text{surrogate}}(y_i, z_i)$$

loss

Test

error

loss