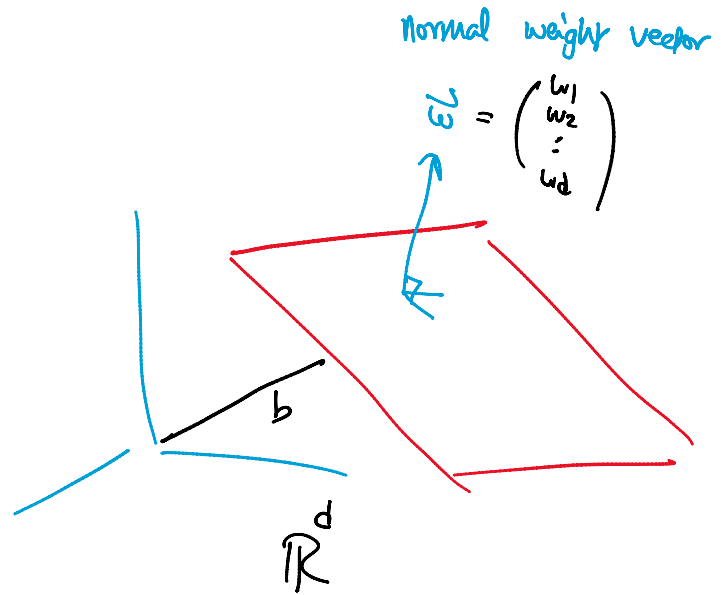
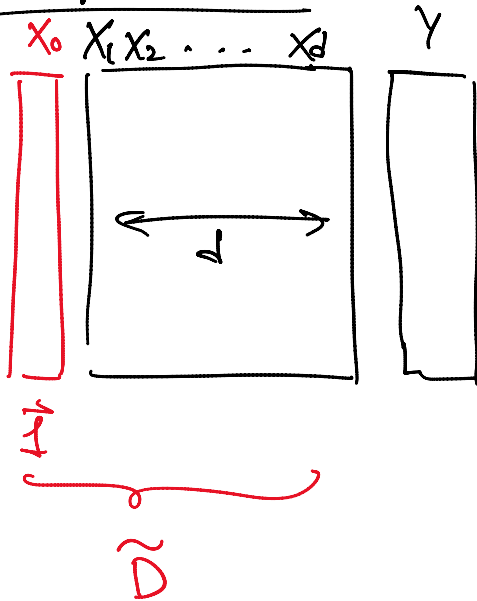


Lecture 10

Thursday, September 28, 2023 10:01 AM

Multiple Regression



$$\hat{y}_i = \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i \quad \forall i = 1, \dots, n$$

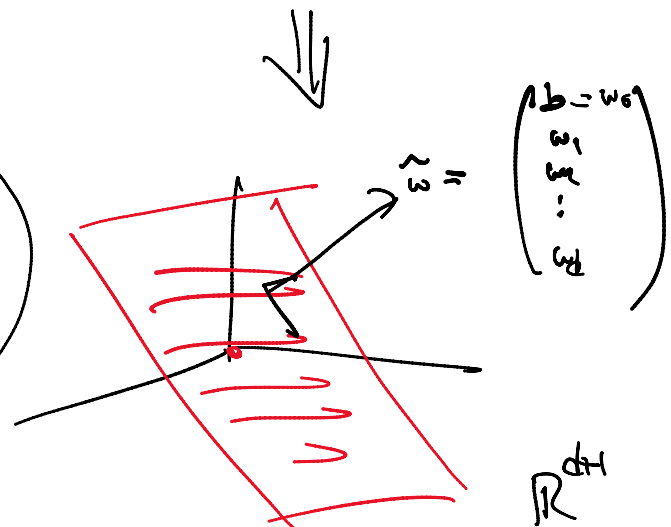
$$= w_0 x_0 + w_1 x_1 + \dots + w_d x_d$$

$$= b + w_1 x_1 + \dots + w_d x_d$$

$$\hat{\mathbf{y}} = \tilde{\mathbf{D}} \tilde{\mathbf{w}}$$

Diagram showing the dimensions of the matrices:

- $\tilde{\mathbf{D}}$ is $n \times (d+1)$.
- $\tilde{\mathbf{w}}$ is $(d+1) \times 1$.



point-view

$$\|\mathbf{y} - \hat{\mathbf{y}}\|^2 = \text{SSE}$$

$$= \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

11.1.1.1.2

$$SSE = \|Y - \hat{Y}\|^2$$

$$= (Y - \hat{Y})^T (Y - \hat{Y})$$

$$= \|Y\|^2 - 2\hat{Y}^T Y + \hat{Y}^T \hat{Y}$$

$$= \|Y\|^2 - 2(\tilde{D}\tilde{\omega})^T Y + (\tilde{D}\tilde{\omega})^T (\tilde{D}\tilde{\omega})$$

$$SSE = \|Y\|^2 - 2\tilde{\omega}^T (\tilde{D}^T Y) + \tilde{\omega}^T (\tilde{D}^T \tilde{D}) \tilde{\omega}$$

$$\left[\frac{\partial SSE}{\partial \tilde{\omega}} = 0 \right]$$

$$\frac{\partial SSE}{\partial \tilde{\omega}} = -2\tilde{D}^T Y + 2\tilde{D}^T \tilde{D} \tilde{\omega} = 0$$

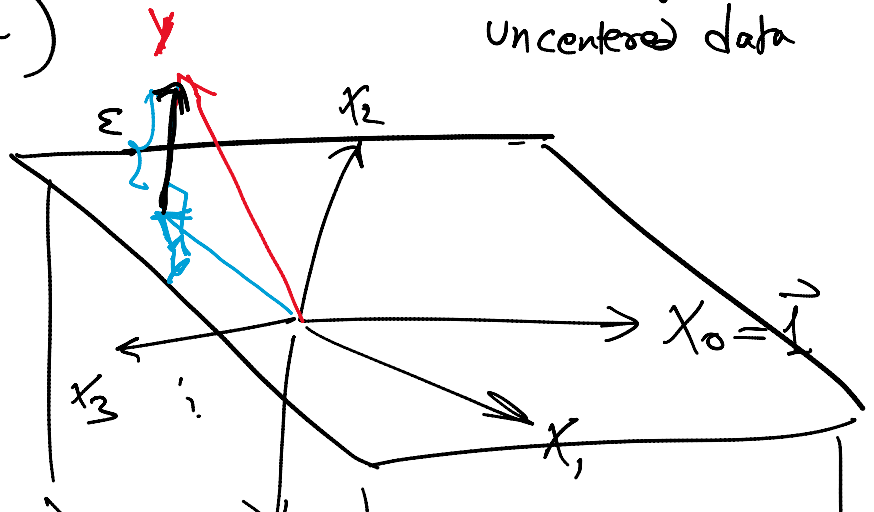
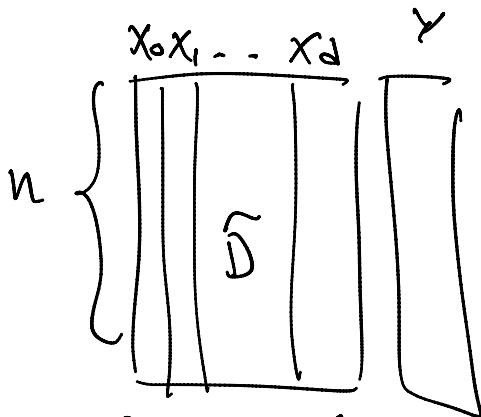
$$\tilde{D}^T \tilde{D} \tilde{\omega} = \tilde{D}^T Y$$

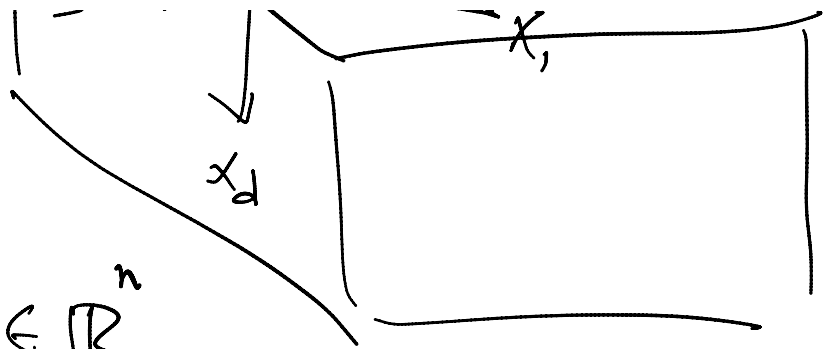
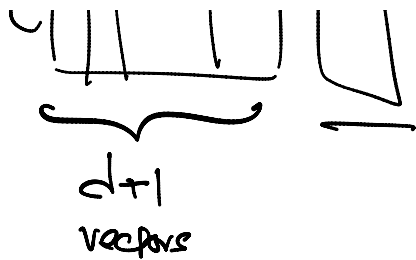
normal equations

$$\tilde{\omega} = (\tilde{D}^T \tilde{D})^{-1} \tilde{D}^T Y$$

like Cov! scatter for uncentered data

Column view (Geometric)



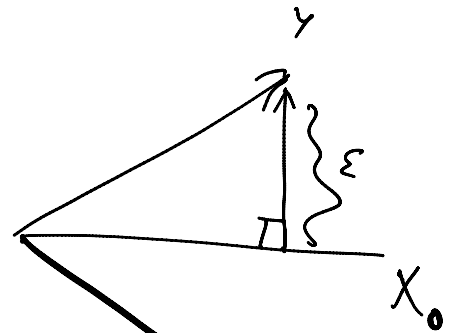


error vector $\vec{\varepsilon} = y - \hat{y} \in \mathbb{R}^n$

$(d+1)$ dim subspace of \mathbb{R}^n

$\vec{\varepsilon}$ is orthogonal to the subspace

$\vec{\varepsilon}$ is orthogonal to all vectors
e.g. x_i



$$\tilde{D}^T \vec{\varepsilon} = 0$$

$$\tilde{D} = \begin{bmatrix} | & | & & | \\ x_0 & x_1 & \dots & x_d \\ | & | & & | \end{bmatrix}$$

$$\tilde{D}^T = \begin{bmatrix} x_0^T & - & - \\ x_1^T & - & - \\ \vdots & & \\ x_d^T & - & - \end{bmatrix}_{d+1 \times n}$$

$$\tilde{D}^T \vec{\varepsilon} = 0$$

$$\Rightarrow \tilde{D}^T (y - \hat{y}) = 0$$

$$\Rightarrow \tilde{D}^T y = \tilde{D}^T \hat{y}$$

$$\tilde{D}^T y = \tilde{D}^T (\tilde{D} \tilde{w})$$

$$\tilde{D}^T \tilde{D} \tilde{w} = \tilde{D}^T y$$

normal equation

$$\begin{bmatrix} x_0^T \vec{\varepsilon} \\ x_1^T \vec{\varepsilon} \\ \vdots \\ x_d^T \vec{\varepsilon} \end{bmatrix} = 0$$

$$\tilde{D}^T \tilde{D} \tilde{w} = \tilde{D}^T y$$

normal equation

$\{X_0 \ X_1 \ \dots \ X_d\} \leftarrow$ not a good basis

\Downarrow

$\{V_0 \ V_1 \ \dots \ V_d\} \leftarrow$ orthogonal basis

$$\hat{y} = \underbrace{\text{proj}_{V_0}(y)} \cdot \vec{V}_0 + \underbrace{\text{proj}_{V_1}(y)} \cdot \vec{V}_1 + \dots + \underbrace{\text{proj}_{V_d}(y)} \cdot \vec{V}_d$$

via QR

$$\tilde{D}^T \tilde{D} \tilde{w} = \tilde{D}^T y$$

$$\tilde{w} = \underbrace{(\tilde{D}^T \tilde{D})^{-1}}_{\substack{\text{may not exist} \\ \text{Positive semi-definite}}} (\tilde{D}^T y)$$

Ridge regression

$$\hat{w} = \left(\tilde{D}^T \tilde{D} + \alpha I \right)^{-1} \tilde{D}^T y$$

$\alpha > 0$

always have an inverse

full rank matrix

$$\begin{array}{c} \begin{array}{|c|c|c|c|c|} \hline & & & & \\ \hline & & & & \\ \hline & & & & \\ \hline & & & & \\ \hline \end{array} & + & \begin{array}{|c|c|c|c|c|} \hline \alpha & & & & \\ & \alpha & & & \\ & & \alpha & & \\ & & & \alpha & \\ & & & & \alpha \\ \hline \end{array} & = & \begin{array}{|c|c|c|c|c|} \hline \alpha & & & & \\ & \alpha & & & \\ & & \alpha & & \\ & & & \alpha & \\ & & & & \alpha \\ \hline \end{array} \\ \tilde{D}^T \tilde{D} & & \alpha I & & (\tilde{D}^T \tilde{D} + \alpha I) \\ & & \underbrace{\hspace{2cm}} & & \lambda_i + \alpha \\ & & \text{ridge} & & \end{array}$$

If λ_i is an eigenvalue of $\tilde{D}^T \tilde{D}$ then $\lambda_i + \alpha$ is the corresponding eigenvalue of $\tilde{D}^T \tilde{D} + \alpha I$

$$\lambda_i = 0 \Rightarrow \lambda_i + \alpha > 0$$

Ridge Regression

(Regularization)

$$SSE = \|y - \hat{y}\|^2$$

$$\min_{\tilde{w}} J = \|y - \tilde{D} \tilde{w}\|^2 \quad \text{original objective}$$

$$\min_{\tilde{\omega}} J = \|\gamma - \tilde{D} \tilde{\omega}\| \quad \text{original objective}$$

On the training sample

$$\hat{y}_i = \frac{w_0}{s} x_0 + \tilde{w}_1 x_1 + \dots + \frac{w_d}{s} x_d$$

\parallel
loss

$$\min_{\tilde{\omega}} J = \|\gamma - \tilde{D} \tilde{\omega}\|^2 + \alpha \|\tilde{\omega}\|^2$$

$\alpha \equiv$ regularization constant $\alpha \geq 0$

$$\frac{\partial J}{\partial \tilde{\omega}} = -2 \tilde{D}^T \gamma + 2 \tilde{D}^T \tilde{D} \tilde{\omega} + 2 \alpha \tilde{\omega} = 0$$

$$\Rightarrow \underline{\tilde{D}^T \tilde{D}} \tilde{\omega} + \underline{\alpha} \tilde{\omega} = \tilde{D}^T \gamma$$

$$(\tilde{D}^T \tilde{D} + \alpha \underline{I}) \tilde{\omega} = \tilde{D}^T \gamma$$

$$\tilde{\omega} = \underline{(\tilde{D}^T \tilde{D} + \alpha \underline{I})^{-1}} \tilde{D}^T \gamma$$

ridge matrix

Exam I

In person / in class

In person / in class
paper - pencil
Open notes

Chapter 1: dot product
norms
projections
prob concepts

Chapter 2

Mean

Variance

Covariance matrix

eigenvalues / eigenvectors | power iteration

Col / point view

geometric interpretation of μ , Σ , σ , ρ (correlation)

total variance & relation to σ_i^2 & change of basis
& eigenvalues

Normal Distribution

$\rightarrow \Sigma, \Sigma^{-1}$, basis vectors, PCA

Chapter 7: PCA (projected variance/error)

7.1 & 7.2

max variance \equiv min MSE

eigen decomposition

selecting best subspace

$$\sum u = 14$$

(no kernel, svd)

Chapter 20 : LDA (projected means & variances)

Sec 20.1

Optimization objective $\frac{(m_1 - m_2)^2}{s_1^2 + s_2^2}$

$$\Rightarrow \boxed{Bw = \lambda Sw} \rightarrow \text{eq. 20.10}$$

\vec{w} : what it means.

$$\vec{w} = \vec{S}^{-1} (m_1 - m_2)$$

Chapter 6 : high dimensional data

hyper - cubes, spheres, planes

Q! boundary, corners, diagonals, normal distribution

$d=2$

$d=3$

finally $d \rightarrow \infty$?

$$\lim_{d \rightarrow \infty}$$

Chapter 23 : linear regression

23.1, 23.2, 23.3, 23.4

bivariate
multiple

$$QR \leftarrow \text{decomposition}$$
$$\vec{w} : \text{back solve}$$

adding ridge

SVD : gradient

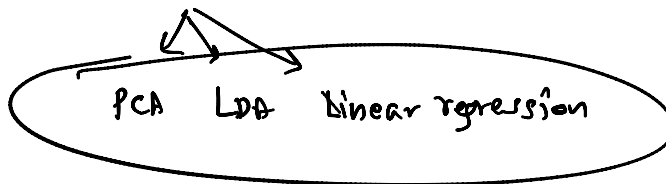


$\vec{w}_0 \equiv \text{initial vector}$

SGD: gradient

→ mean

→ $\vec{w} \leftarrow$ via SGD



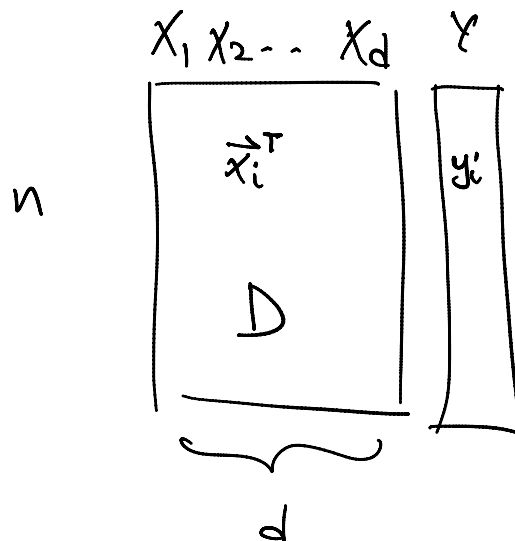
$\hat{w}_0 \equiv$ initial vector

∇_w, η

$$\hat{w}_t = \hat{w}_{t-1} - \eta \cdot \nabla_{w_{t-1}}$$

Logistic Regression

→ binary logistic regression



$y_i \in \{P, N\}$
binary 'class' ↓

1	0
---	---

encoding
+1 -1

$$D = \left\{ (\vec{x}_i^T, y_i) \right\}_{i=1}^n$$

$y_i \in \{0, 1\}$

Q: Should we be using numeric predictions?

A: Not really.

$y_i = 1$, prediction $\hat{y}_i = 1000$

$$(y_i - \hat{y}_i)^2 = \underline{\underline{(\eta_i)^2}}$$

2 classes

$$P(1 | \vec{x}_i) \quad \text{vs} \quad P(0 | \vec{x}_i)$$

given point/instance \vec{x}_i

Task: predict the probability of each class given \vec{x}_i

$$P(1 | \vec{x}_i) + P(0 | \vec{x}_i) = 1$$

$$P(0 | \vec{x}_i) = 1 - P(1 | \vec{x}_i)$$

$$\text{given } \vec{x}_i, \text{ predict } P(1 | \vec{x}_i) \equiv \pi(x_i)$$

prob of positive class
Unknown!

$$P(1 | \vec{x}_i) = \pi(x_i) \neq (\underbrace{w_0}_{\text{Prob } [0,1]}) + w_1 x_{i1} + w_2 x_{i2} + \dots + w_d x_{id}$$

$\underbrace{\vec{w}^T \vec{x}_i}_{\text{dot product}}$

$$[0, 1]$$

$$[-\infty, \infty]$$

$\Downarrow \theta \leftarrow$ Some link function

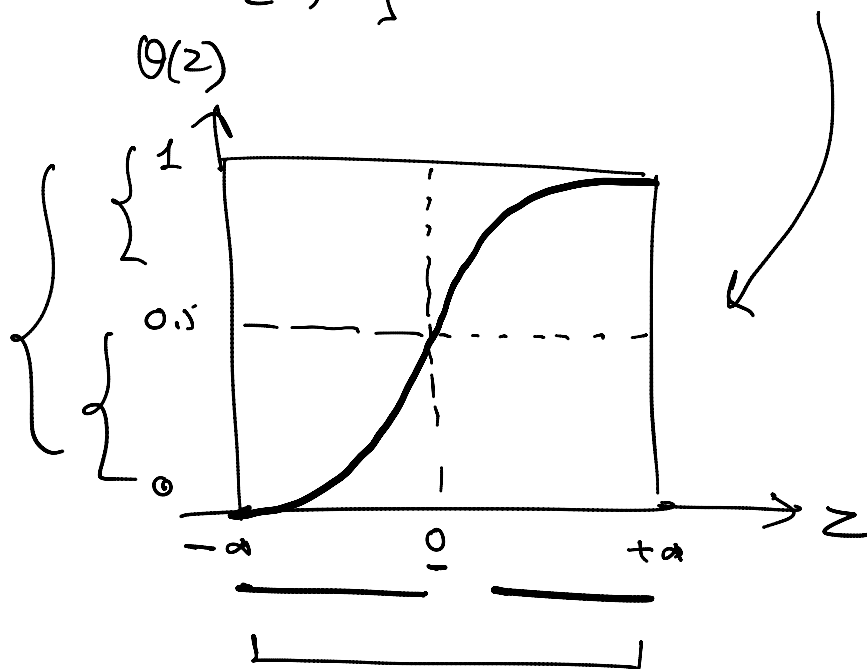
$\theta \equiv$ logistic / sigmoid

$$[0, 1]$$

$$\theta: [-\infty, \infty] \Rightarrow [0, 1]$$

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

↑
Scalar



$$\pi(\vec{x}_i) = \theta(\underbrace{\tilde{\omega} \vec{x}_i}_{\text{scalar}})$$

\vec{x}_i is one of the input points
parameters: $\tilde{\omega}$ unknown

$$\forall i = 1 \dots n$$

n -equations

non-linear model

Q: how to learn $\tilde{\omega}$?

A: maximum likelihood estimation (MLE)

$$\max_{\tilde{\omega}} L(y|\tilde{\omega}) \equiv \text{likelihood of the data}$$

Choose $\tilde{\omega}$ that maximizes the probability of the data that is observed

$$L(Y|\tilde{\omega}) \equiv P(Y|\tilde{\omega})$$

Prob of observed data given $\tilde{\omega}$

$Y \in \mathbb{R}^n$ n -dim joint probability

$$\equiv P(y_1 y_2 \dots y_n | \tilde{\omega}) ?$$

lets make some assumptions!

Assume all points are independent

$$\equiv P(y_1|\tilde{\omega}) \times P(y_2|\tilde{\omega}) \times \dots \times P(y_n|\tilde{\omega})$$

$$L(\tilde{\omega}) = \prod_{i=1}^n P(y_i|\tilde{\omega})$$

$$P(y_i=1|\tilde{\omega}) = \Theta(\tilde{\omega}^T \tilde{x}_i)$$

Case:

$$P(y_i|\tilde{\omega}) = \begin{cases} \Theta(\tilde{\omega}^T \tilde{x}_i) & \text{if } y_i=1 \\ 1 - \Theta(\tilde{\omega}^T \tilde{x}_i) & \text{if } y_i=0 \end{cases}$$

$$P(y_i | \tilde{w}) = \theta(\tilde{w}^T x_i)^{y_i} \cdot (1 - \theta(\tilde{w}^T x_i))^{1-y_i}$$

$$y_i = 0 \leftarrow$$

$$y_i = 1$$

we
 \mathcal{L}

$$L(\tilde{w}) = \prod_{i=1}^n \theta(\tilde{w}^T x_i)^{y_i} \cdot (1 - \theta(\tilde{w}^T x_i))^{1-y_i}$$

$$\frac{\partial L(\tilde{w})}{\partial \tilde{w}} = 0$$

Solve!