


Lecture 3

DATA 558

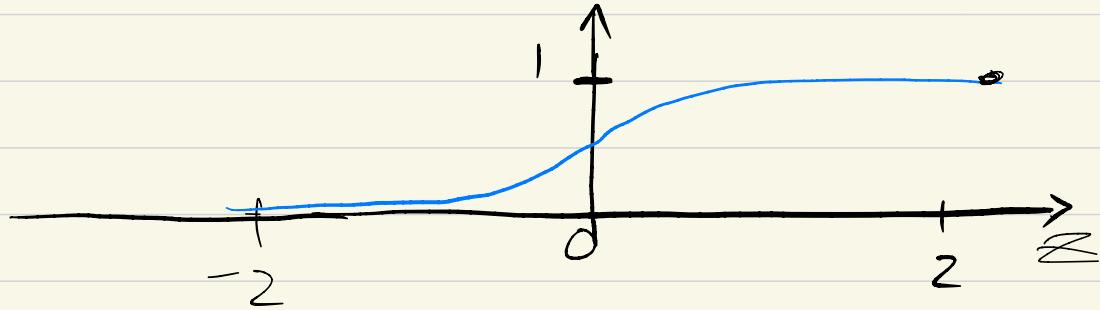
Min
 $\beta \in \mathbb{R}^d$

$$\frac{1}{n} \sum_{i=1}^n \log \left(1 + e^{-y_i \times_i^\top \beta} \right) + \lambda \|\beta\|_2$$

$$P(Y=1 | x; \beta) = g(x^\top \beta)$$

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

g is called logistic function
sigmoid STAT
ML



1

2020

$$\begin{aligned} & \text{IP}(Y=+1 | x; \beta) \\ + & \text{IP}(Y=-1 | x; \beta) \\ = & 1 \end{aligned}$$

exclusive
class
membership

$$\begin{aligned} & \text{IP}(Y=-1 | x; \beta) \\ = & 1 - \text{IP}(Y=+1 | x; \beta) \end{aligned}$$

$$\log \left(\frac{\text{IP}(Y=+1 | x; \beta)}{\text{IP}(Y=-1 | x; \beta)} \right)$$

$$= x^T \beta$$

$g(z)$ nice property

2

$$1 - g(z) = g(-z)$$

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

$$= \frac{1}{1 + e^{+z}}$$

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

$$e^{-z}$$

$$= \frac{e^{-z}}{e^{-z} + e^{\cancel{-z} + z}}$$

(3)

$$1 - g(+z) = g(-z)$$

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

$$\begin{aligned} & e^{-z} + 1 - 1 \\ &= \frac{e^{-z}}{1 + e^{-z}} \\ &= \frac{\cancel{e^{-z} + 1}}{\cancel{1 + e^{-z}}} - \frac{1}{1 + e^{-z}} \\ &= 1 - \frac{1}{1 + e^{-z}} \end{aligned}$$

$$g(-z) = 1 - g(+z)$$

(4)

$$\text{IP} (Y=+1 \mid x; \beta)$$
$$= g(z) = g(x^T \beta)$$

$$\text{IP} (Y=-1 \mid x; \beta)$$
$$= 1 - g(z)$$

$$= g(-z) = g(-x^T \beta)$$

(S)

Un-regularized Maximum Likelihood

Estimator for Logistic Regression

Likelihood \mathcal{L}

$$\mathcal{L}((x_1, y_1), \dots, (x_n, y_n))$$

$$= \prod_{i=1}^n p(y_i | x_i; \beta)$$

↑

$$p(y_1 | x_1; \beta) \times \dots \times p(y_n | x_n; \beta)$$

$$\underset{\beta \in \mathbb{R}^d}{\operatorname{Max}} \quad \mathcal{L}(\beta)$$

$$\underset{\beta \in \mathbb{R}^{d+1}}{\operatorname{Min}} -\log(\mathcal{L}(\beta)).$$

6

$$\begin{aligned}
 & -\log \left(\prod_{i=1}^n p(y_i | x_i, \beta) \right) \\
 & = -\sum_{i=1}^n \log (p(y_i | x_i; \beta))
 \end{aligned}$$

$$\log(a \times b) = \log(a) + \log(b)$$

$$\begin{aligned}
 & -\log \left(\prod_{i=1}^n \dots \right) \\
 & = \sum_{i=1}^n -\log p(y_i | x_i; \beta) \\
 & = \sum_{i=1}^n -\log g(y_i \cdot x_i^T \beta)
 \end{aligned}$$

$$\begin{aligned}
 p(y=+1 | x_i, \beta) & = g(x_i^T \beta) = g(+1 \cdot x_i^T \beta) \\
 p(y=-1 | x_i; \beta) & = 1 - g(x_i^T \beta) \\
 & = g(-x_i^T \beta) = g(-1 \cdot x_i^T \beta)
 \end{aligned}$$

(7)

$$\text{IP}(Y=y | x_i \beta) = g(Y_i \cdot x_i^T \beta)$$

where $y = +1$ or $y = -1$

- log (Likelihood)

$$= \sum_{i=1}^n -\log(g(Y_i \cdot x_i^T \beta)) \quad \begin{array}{l} \text{(all terms} \\ \text{write that} \\ \text{way)} \end{array}$$

$$= \sum_{i=1}^n -\log \left(\frac{1}{1 + e^{-Y_i \cdot x_i^T \beta}} \right)$$

$$\text{but } \log\left(\frac{1}{a}\right) = -\log(a)$$

$$= \sum_{i=1}^n (-1) \times (-1) \cdot \log\left(\frac{1}{1 + e^{-Y_i \cdot x_i^T \beta}}\right)$$

$$\frac{1}{n} \sum_{i=1}^n \log\left(\frac{1}{1 + e^{-Y_i \cdot x_i^T \beta}}\right)$$

Empirical risk / loss

(8)

$$\begin{aligned}
 & \text{Min}_{\beta \in \mathbb{R}^d} f(\beta) \\
 &= \left[\frac{1}{n} \sum_{i=1}^n \log \left(1 + e^{-y_i x_i^T \beta} \right) \right. \\
 &\quad \left. + \lambda \|\beta\|_2^2 \right]
 \end{aligned}$$

Gradient descent algorithm

$$\beta_{t+1} = \beta_t - \eta_t \nabla_{\beta_t} F$$

eta

3

Single term $\frac{\partial}{\partial \beta} \log(1 + e^{-y \cdot x^T \beta})$

$$l(\cdot) = \log(\cdot)$$

$$g(\cdot) = 1 + \exp(-\cdot)$$

$$h(\cdot) = y \cdot x^T \beta$$

$$l(g(h(\beta)))$$

$$l((g \circ h)(\beta))$$

$$(l \circ g \circ h)' = (g \circ h)' \cdot l'((g \circ h))$$

$$(g \circ h)' = h' \cdot g'(h)$$

(10)

$$(f \circ g \circ h)' =$$

$$(1 + \exp(-bx))^c = -\exp(-bx)$$

$$= (h') \cdot g'(h) \cdot f'(g \circ h)$$

$$\boxed{y \cdot x}$$

$$-1 \cdot e$$

$$-(y \cdot x)^T \beta$$

$$\frac{\partial}{\partial b} (\sqrt{b}x) = v$$

$$\frac{1}{1 + e^{-y \cdot x^T \beta}}$$

$$\frac{\partial}{\partial \beta} \log \left(\frac{1}{1 + e^{-y \cdot x^T \beta}} \right)$$

$$= -y \cdot x \cdot \frac{e^{-y \cdot x^T \beta}}{1 + e^{-y \cdot x^T \beta}}$$

11

$$\frac{\partial}{\partial \beta} \left\{ \log \left(\frac{e^{-y_i x_i^T \beta}}{1 + e^{-y_i x_i^T \beta}} \right) \right\}$$

$$= -y_i \cdot x_i \cdot \frac{e^{-y_i x_i^T \beta}}{1 + e^{-y_i x_i^T \beta}}$$

$$1 - p_i = \frac{e^{-y_i x_i^T \beta}}{1 + e^{-y_i x_i^T \beta}}$$

$$\frac{\partial}{\partial \beta} \left\{ \dots \right\} = -y_i x_i \cdot (1 - p_i)$$

$\log - \text{sum} - \exp$

(12)

$$\frac{1}{n} \sum_{i=1}^n \log \left(1 + e^{-y_i x_i^T \beta} \right)$$

$$= \frac{1}{n} ? \bar{Y} \bar{X} ? \circ$$

$$\frac{1}{n} \sum_{i=1}^n (y_i - x_i^T \beta)^2$$

$$= \frac{1}{n} \| \bar{Y} - \bar{X} \beta \|_2^2$$

(13)

$$Q = \text{Id} - \text{Diag}([P_1, \dots, P_n])$$

P_i as defined earlier

$$Q = \begin{pmatrix} 1-P_1 & & & \\ & \ddots & & \\ & & 1-P_n & \\ & & & \ddots \end{pmatrix}$$

$$\underset{\beta \in \mathbb{R}^d}{\text{Min}} \quad \frac{1}{n} \sum_{i=1}^n \log(1 + e^{-y_i x_i^T \beta}) R(\beta)$$

$$+ \lambda \|\beta\|_2^2$$

$$\frac{\partial}{\partial \beta} R(\beta) = \frac{1}{n} \sum_{i=1}^n -y_i x_i (1 - p_i)$$

$$= -\frac{1}{n} Q Y X$$

(14)

$$\frac{\partial}{\partial \beta} \left\{ \lambda \|\beta\|_2^2 \right\} = ?$$

$$\frac{\partial F}{\partial \beta} = \boxed{-\frac{1}{n} Q Y X + 2 \lambda \beta}$$

—
Regularized logistic regression

Penalized Maximum Likelihood predictor

Mathematical derivation

Gradient Computation

Smoothness constant L

$$\|\nabla F(\beta) - \nabla F(\beta')\|_2^2 \leq L \|\beta - \beta'\|_2$$

For all $\beta, \beta' \in \mathbb{R}^d$

Gross estimate of L

Pick two β, β' at random.

$$L_{\text{est}} \approx \frac{\|\nabla F(\beta) - \nabla F(\beta')\|_2^2}{\|\beta - \beta'\|_2}$$

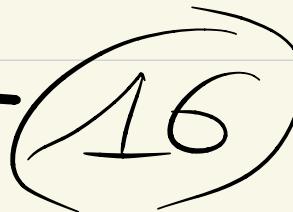
quick

Better estimate of L

$L_{\text{est}} \approx$ maximum eigenvalue of

$$\left\{ \left(\frac{1}{m} X_{\text{sub}}^T X_{\text{sub}} \right) + \lambda \text{Id} \right\}$$

Python largest-eigh

slow 

Adaptive step-size rule

Back-tracking line-search

Init $\gamma_t = \gamma_{t-1}$, if $t \geq 1$

$$\gamma_0 = \frac{1}{L_{\text{est}}}$$

Iterate

$$\gamma_t = \gamma \cdot \gamma_t$$

until condition is satisfied

$$F(\beta_{t+1}) - F((\beta_t) + \gamma_t \nabla F(\beta_t))$$

$$\leq F(\beta_t) - \alpha \gamma_t \|\nabla F(\beta_t)\|^2$$

where $\alpha = 0.5$ and $\gamma = 0.8$

First Gradient Method

Init $\beta_0 = \theta$; $\theta_0 = \theta$

Iterate Find γ_t by backtracking
line-search

$$\beta_{t+1} = \theta_t - \gamma_t \nabla_{\theta} F(\theta)$$

Momentum $\theta_{t+1} = \beta_{t+1} + \frac{\epsilon}{t+3} (\beta_{t+1} - \beta_t)$

$$(\beta_t)_{t \geq 1}$$

$$(\theta_t)_{t \geq 1}$$

two sequences
of iterates
at the same time

$$O\left(\frac{1}{\epsilon}\right) \longrightarrow O\left(\frac{1}{\epsilon^2}\right)$$

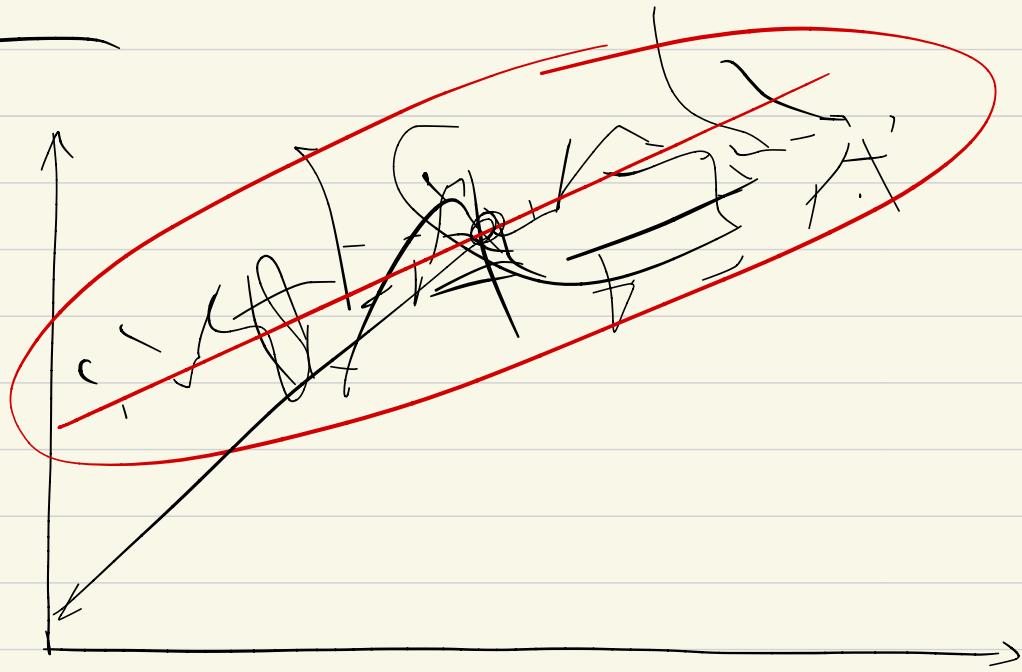
(18)

Nesterov accelerated gradient method

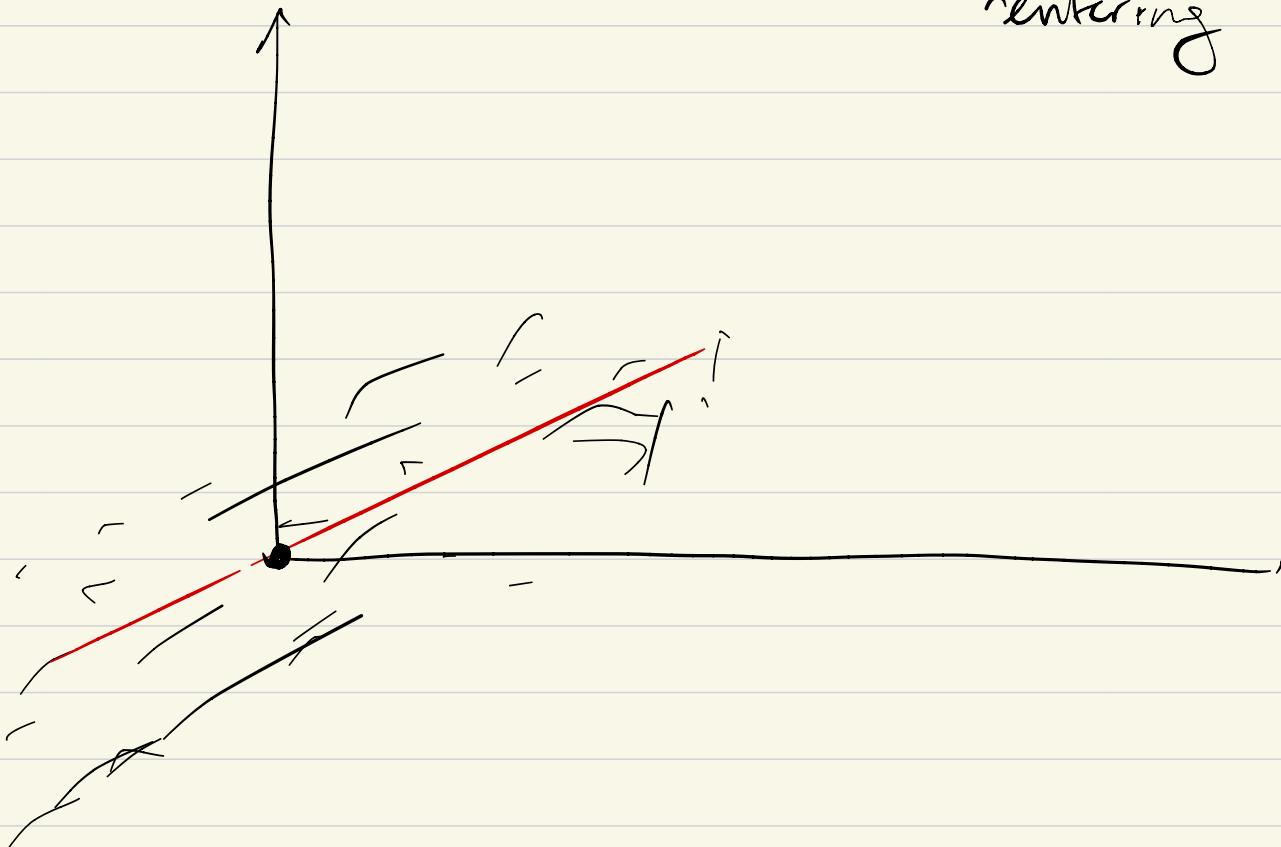
Polyak momentum



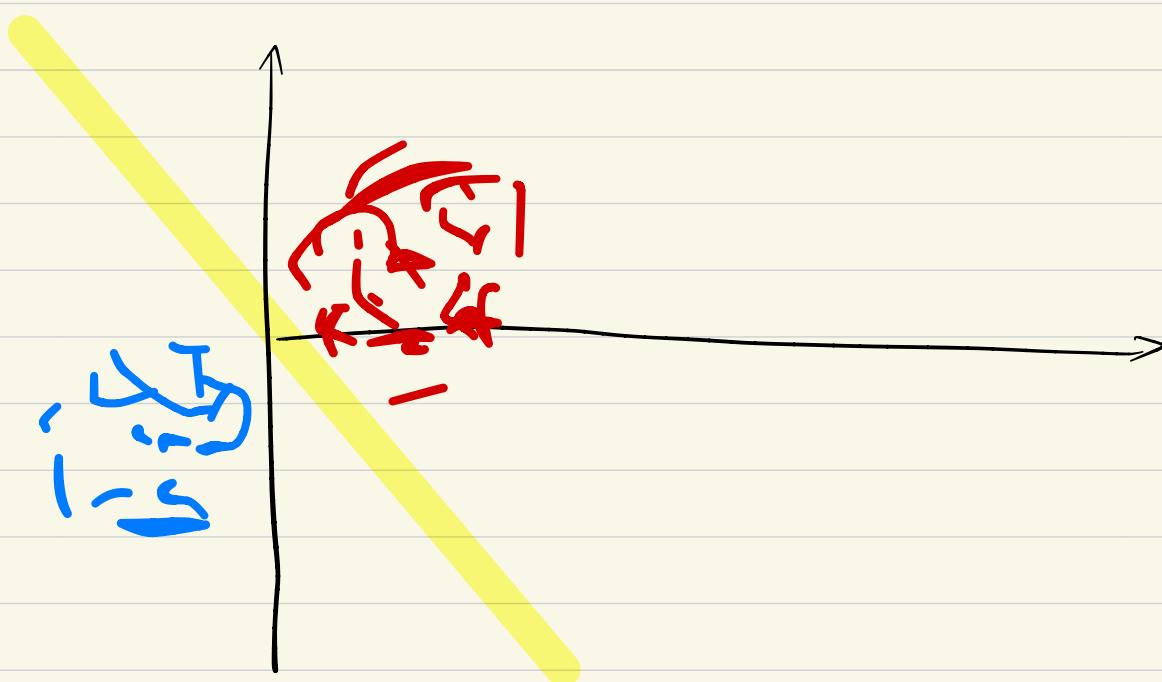
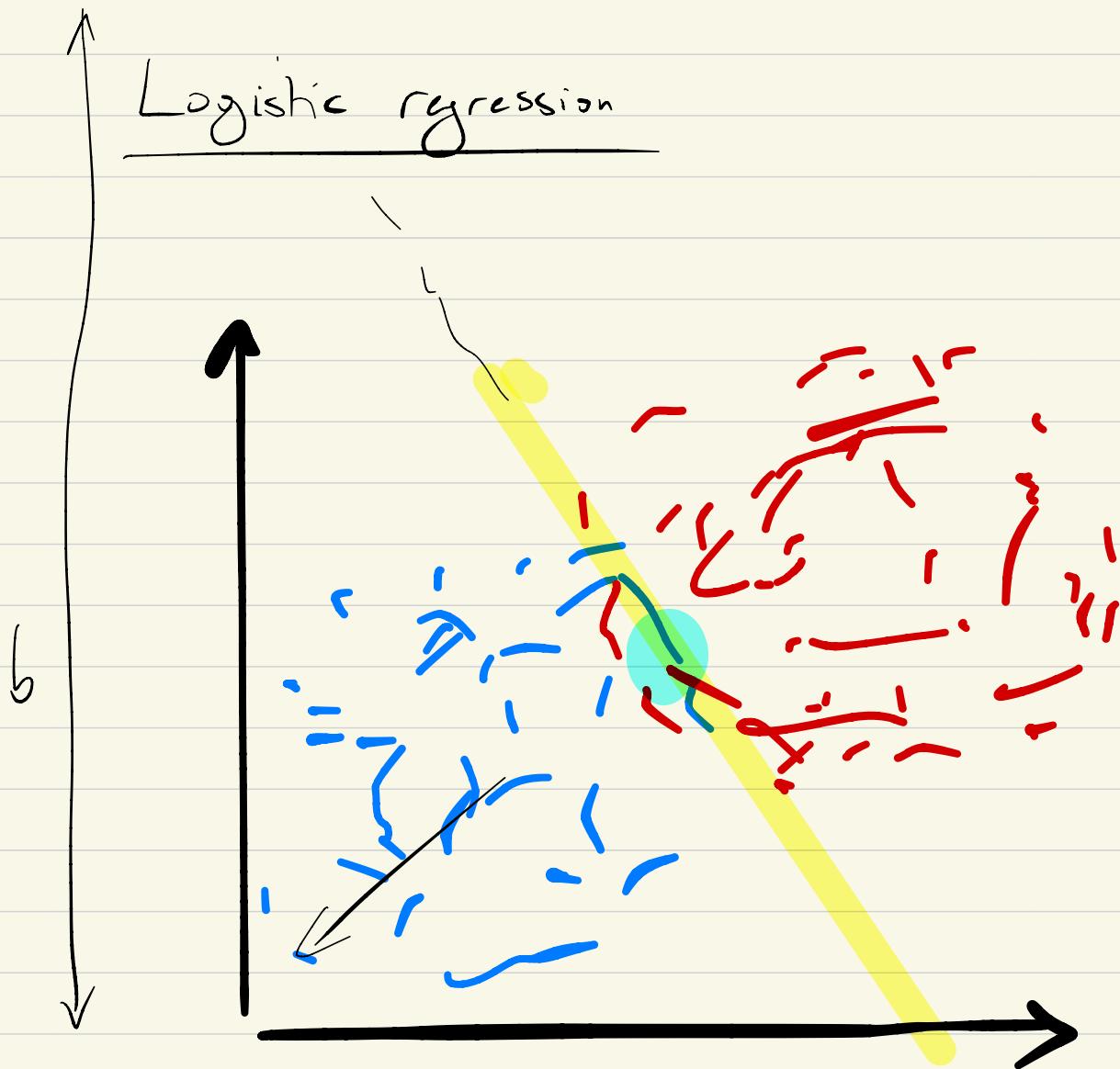
Regression



centering



Logistic regression



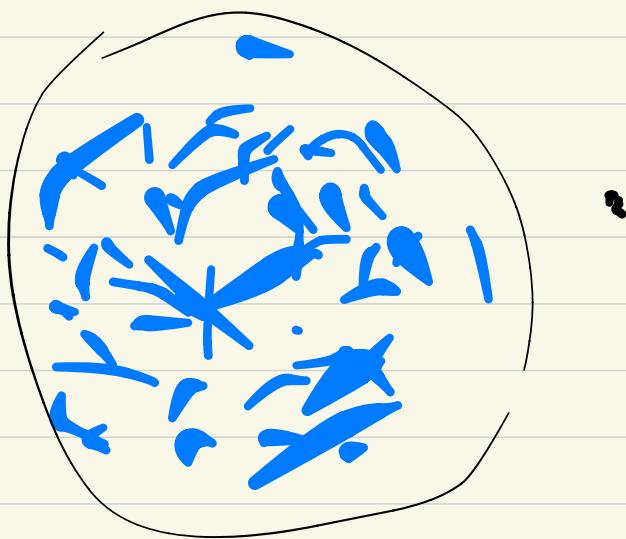
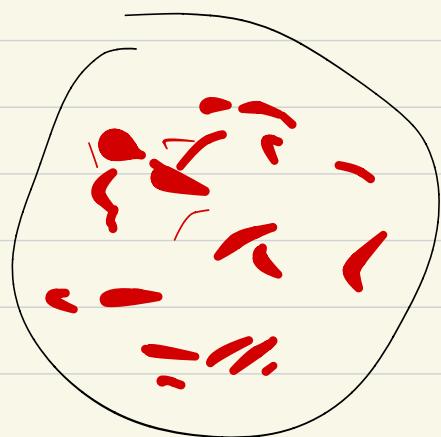
Training data

non-Spam



spam

Testing data



$$\tilde{x} = \begin{bmatrix} \vdots \\ 1 \end{bmatrix} \quad \tilde{\beta} = \begin{bmatrix} \vdots \\ \downarrow \\ \text{Intercept} \end{bmatrix}$$

$$\boxed{\tilde{x}^T \tilde{\beta}} = x^T \beta$$

actual features

+ 1 x intercept

$$= x^T \beta + \text{intercept}$$

$$\|\tilde{\beta}\|_2^2 = \sum_{j=1}^d \tilde{\beta}_j^2 + \text{intercept}^2$$