

Today: Start classification

Recap regression via least squares methods

Find $f: \mathbb{R}^d \rightarrow \mathbb{R}$ so that

$$f(\vec{x}_i) \approx y_i \quad \text{for } (\vec{x}_i, y_i)_{i=1}^n$$

$$\min_f \frac{1}{n} \sum_{i=1}^n (f(\vec{x}_i) - y_i)^2$$

training data

Mean squared error \rightarrow least squares estimator

Models for f

- linear model $f(\vec{x}) = \beta_0 + \beta_1 x_1 + \dots + \beta_d x_d$

intercept
offset

coefficient

$$= [1, x_1, \dots, x_d] \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_d \end{bmatrix}$$

$$= \vec{x}^T \vec{\beta}$$

- fixed features + linear model

$$f(\vec{x}) = \vec{\Phi}(\vec{x})^T \vec{\beta} = \sum_{i=1}^f \phi_i(\vec{x}) \beta_i$$

basis



ex/ degree-2 polynomials

basis fns

$$\vec{\Phi}(\vec{x}) = [1, x_1, x_2, x_1^2, x_1 x_2, x_2^2]^T \in \mathbb{R}^6$$

$$\vec{x} \in \mathbb{R}^2, d=2, \vec{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

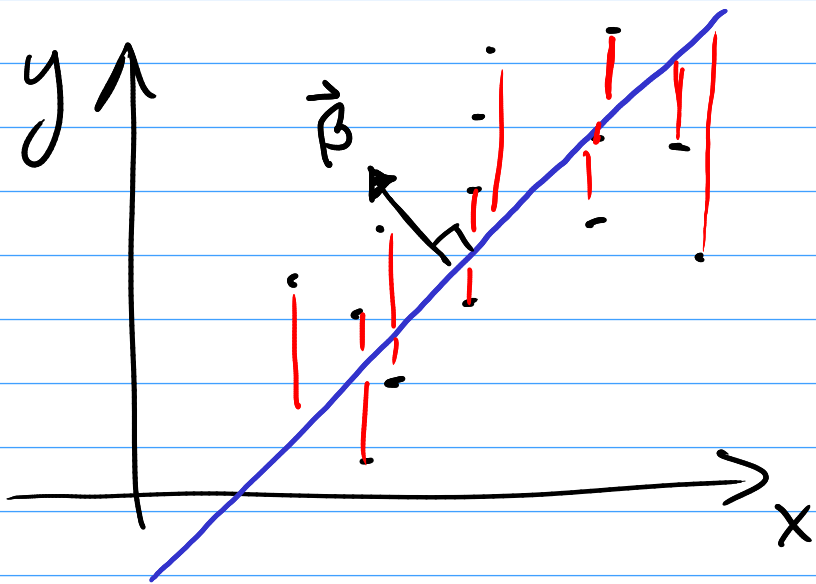
parametric
model
= fixed
dimensions $\vec{\beta}$

feature
vector

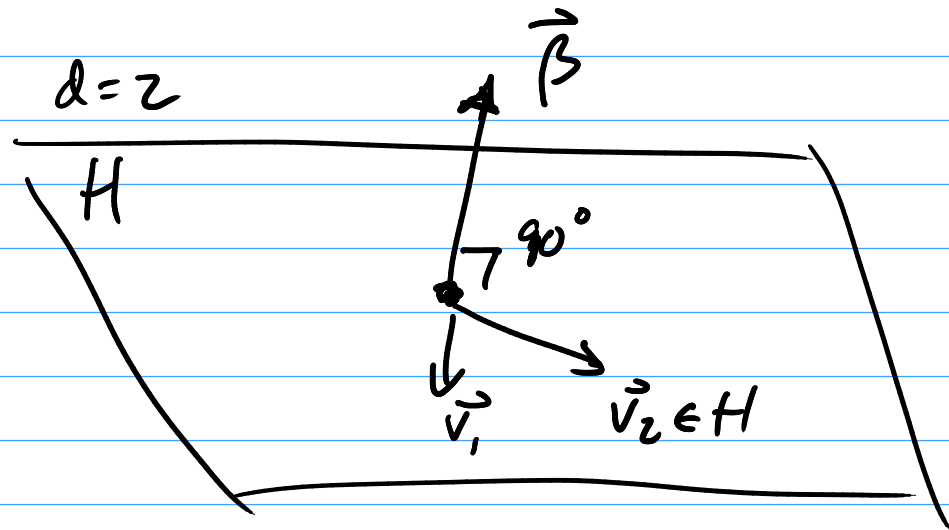
featurization

Linear predictors fit hyperplane to (\vec{x}_i, y)

$d=1$



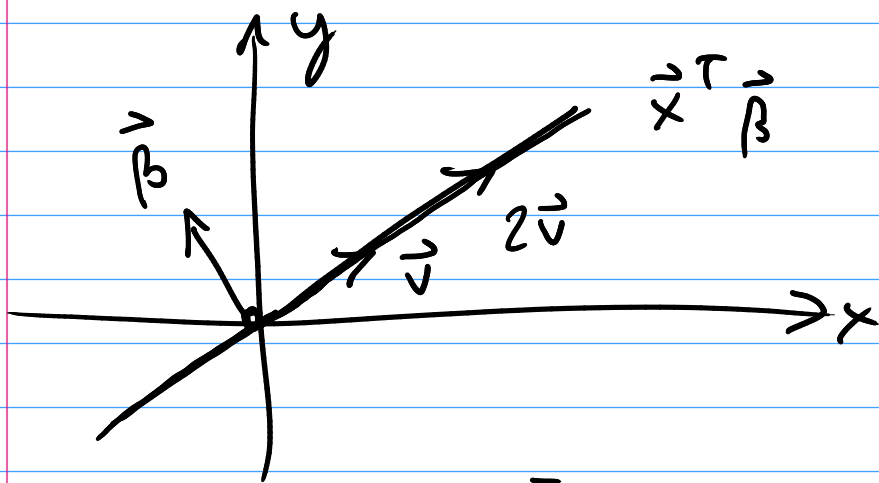
in $d=1$ dims, hyperplane
is a line



$\vec{\beta}$ is orthogonal
to the hyperplane

$$\vec{\beta}^T \vec{v} = 0 \text{ for any } \vec{v} \in H$$

If $\beta_0 = 0$ then $\vec{\beta}$

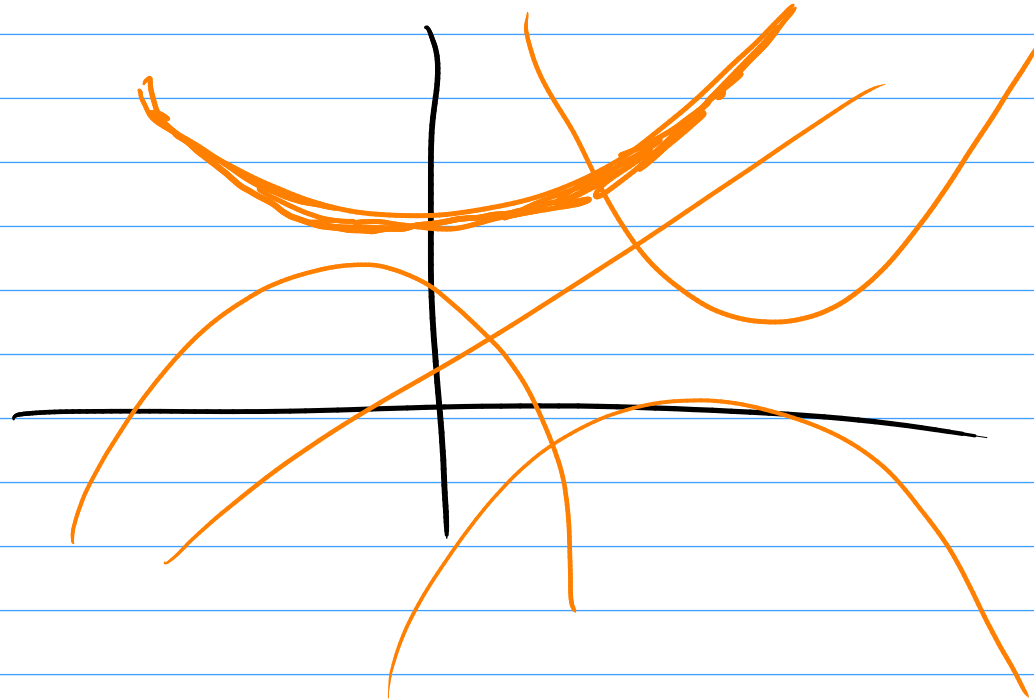


orthogonal: $\vec{\beta}^T \vec{v} = 0$

orthonormal: $\|\vec{\beta}\| = 1$
 $\|\vec{v}\| = 1$

$\phi_i(\bar{x})$

$$3 \cdot 1 - \frac{1}{2} \cdot x + x^2 = 3 + x^2 - \frac{1}{2}x$$



basis of nonlinear features

Estimators:

- Ordinary Least Squares

$$\min_{\vec{\beta}} \left\| X \vec{\beta} - \vec{y} \right\|^2$$

features labels
↓ ↓

- Ridge regression

controlling variance
noise in y 's
collinear features

- small singular values of X

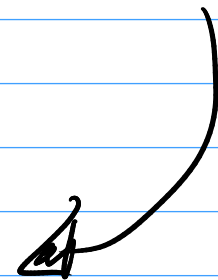
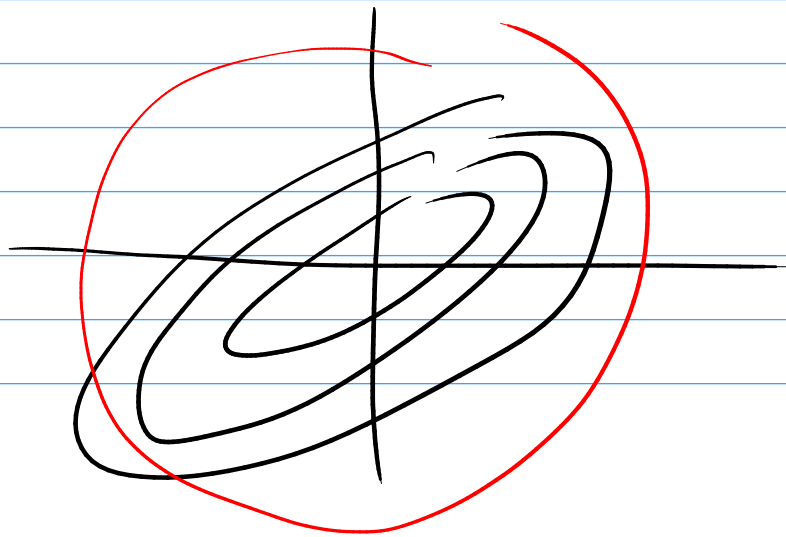
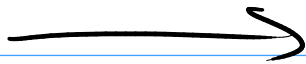
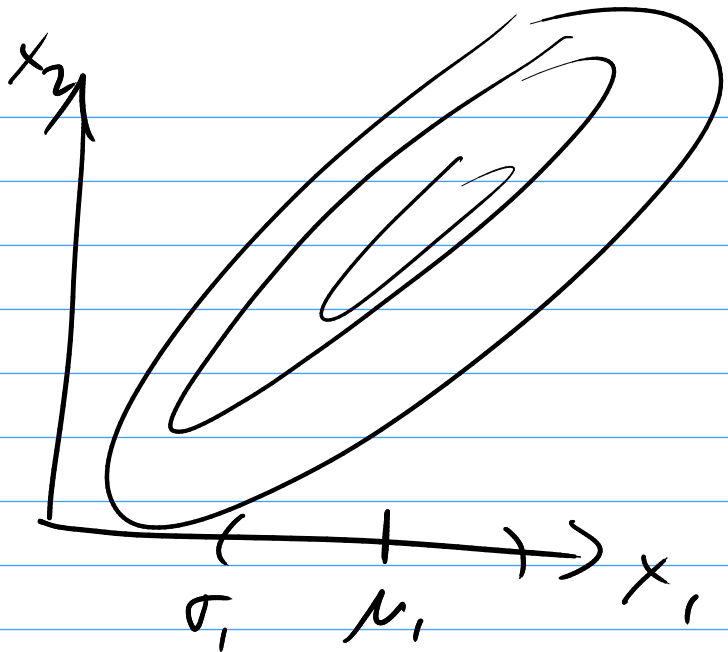
$$\dots + \underbrace{\lambda \|\vec{\beta}\|^2}_{\text{regularization penalty}}$$

- Lasso, giving sparsity

$$\dots + \lambda \|\vec{\beta}\|_1$$

$$\vec{\beta}_{\text{lasso}} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

Preprocessing



Standardization

vs. whitening
makes singular values
equal

Turns out that least squares

Best possible predictor

$$f(\vec{x}) = E[Y | \vec{x}]$$

$$\min_{\hat{y}} (\hat{y} - y)$$

↑
prediction

Bayes predictor

Classification

$$y = \begin{cases} +1 \\ -1 \end{cases}$$

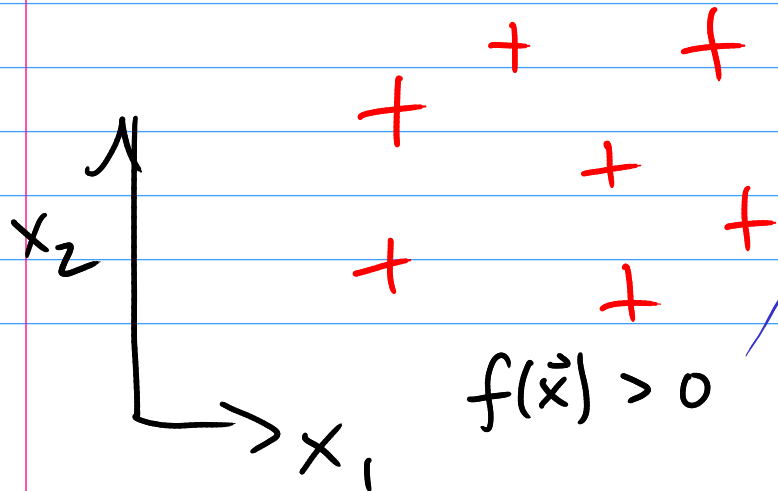
"true"

"false"

categorical

Want f to model

$$\Pr[Y=+1 | \vec{x}]$$

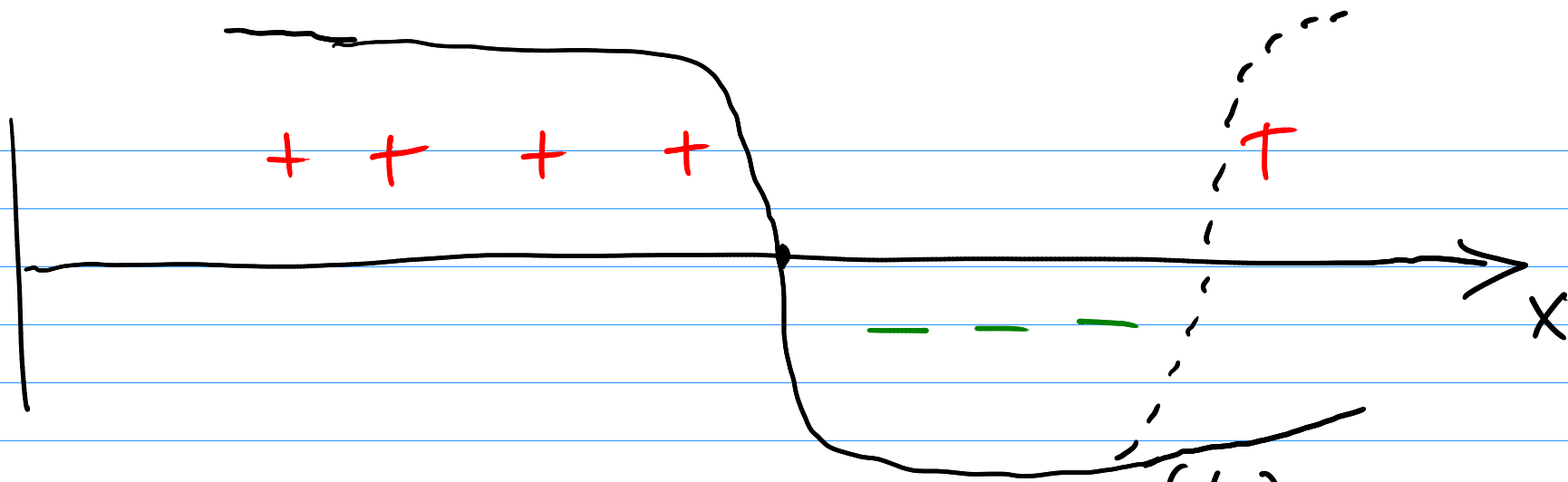


$$f(\vec{x}) = 0$$

$$f(\vec{x}) > 0$$

decision boundary

$$f(\vec{x}) < 0$$



$$f(x_1) = 10$$

$$f(x_2) = 1$$

$$f(x_3) = 0.1$$

Confidence that $y = +1$
 high
 med
 low

5

$$\frac{1}{n} \sum_{i=1}^n (f(\vec{x}_i) - y_i)^2$$

$n \swarrow$ n_{test} vs. n_{train}

$$\frac{1}{n} \|X\hat{\beta} - \vec{y}\|^2$$

arrays
↓ ↓
np.linalg.mean $((y_{\text{pred}} - y) \times 2)$

MSE mean square error

2.3-4 very similar

4.5

$$Y = [2 \times 61]$$

