

Linear regression: minimizing sum of squares of errors in

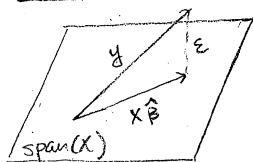
$$\underbrace{\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}}_y = \underbrace{\begin{pmatrix} 1 & x_{11} & x_{12} & \dots & x_{1k} \\ 1 & x_{21} & x_{22} & \dots & x_{2k} \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nk} \end{pmatrix}}_X \underbrace{\begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{pmatrix}}_\beta + \underbrace{\begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}}_\varepsilon$$

i.e. find β s.t. $\|y - X\beta\|$ is minimized
 solution by differentiating wrt β :

$$\hat{\beta} = (X^T X)^{-1} X^T y \quad \text{OLS estimator for } \beta$$

$X^T X$ invertible iff X has linearly indep. cols
 fails when (1) more variables than samples
 (2) perfect multicollinearity

geometric interpretation of OLS



$X\hat{\beta}$ is a linear combination of cols of X

if $y \in \text{span}(X)$, $\varepsilon = 0$

otherwise, ε is minimized when $X\hat{\beta} \perp \varepsilon$, i.e. $X\hat{\beta}$ is the orthogonal projection of y onto $\text{span}(X)$

Projection matrix: a square matrix $P: W \rightarrow W$ s.t. $P^2 = P$ (think geometrically)

in particular, $\forall v \in W, P(v - Pv) = 0$ $W = \text{Im}(P) \oplus \text{Null}(P)$ unique decomposition
 "ker(P)"

orthogonal projection: $\text{Im}(P)$ and $\text{Null}(P)$ are orthogonal subspaces $\Leftrightarrow P$ is self-adjoint, $P = \overline{P^T}$

in particular, $Pv \perp v - Pv \quad \forall v \in W$

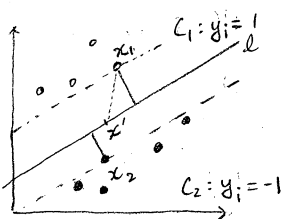
$$P = P^T$$

solving for the orthogonal projection $P: y \rightarrow \text{span}(X)$

$$Py = X\hat{\beta} \text{ for some } \hat{\beta} \text{ and } y - Py \perp \text{span}(X) \Rightarrow X^T(y - X\hat{\beta}) = 0 = X^T y - X^T X \hat{\beta} \Rightarrow \hat{\beta} = (X^T X)^{-1} X^T y$$

PCA

SVM find hyperplane that maximizes the separation between the 2 classes



training data

$$\{(x_i, y_i) \mid x_i \in \mathbb{R}^k, y_i \in \{1, -1\}\}$$

(1). Hard margin: hyperplane exists that separates the classes

- hyperplane l represented by $w_1 x_{(1)} + \dots + w_k x_{(k)} + b = 0$, or $w^T x + b = 0$ $\begin{cases} > 0 & x_i \in C_1 \\ < 0 & x_i \in C_2 \end{cases}$
- let $x_1 \in C_1, x_2 \in C_2$ be the points closest to l . wlog, $w^T x_1 + b = 1$ and $w^T x_2 + b = -1$
 $\Rightarrow y_i (w^T x_i + b) \geq 1$

objective is to maximize the margin around l , i.e. $d(x_1, l) + d(x_2, l)$, where $d(x, l)$ is the projection of $x - x'$ onto the normal vector of l , where x' is any pt on l

to see that w is the normal vector: $w^T(x' - x'') = -b + b = 0$ for $x', x'' \in l$

from before, $\text{Proj}_w(v) = w \underbrace{(w^T w)^{-1}}_s w^T v = \frac{v \cdot w}{w \cdot w} w$ "vector projection"

with $\|\text{Proj}_w(v)\| = \frac{v \cdot w}{w \cdot w} \|w\| = \frac{v \cdot w}{\|w\|}$ "scalar projection" (can also be derived from $a \cdot b = |a||b|\cos\theta$)

$$\Rightarrow d(x_i, l) = |\text{Proj}_w(x_i - x')| = w \cdot (x_i - x') / \|w\| = 1 / \|w\| = d(x_2, l)$$

$$\Rightarrow \text{maximize } 2/\|w\|, \text{ or minimize } \|w\|^2, \text{ subject to } y_i(w^T x_i + b) \geq 1$$

(2). Soft margin: no linear separability

(if $\xi_i > 1$, misclass.)

- introduce slack variables $\xi_i \geq 0$ that measure the degree to which x_i is on the wrong side of the margin
- new objective function is to maximize the margin while minimizing the classes overlap:

$$\min_{w, \xi_i} (\frac{1}{2}) \|w\|^2 + (\frac{c}{n}) \sum \xi_i \quad \text{subject to } y_i (w^T x_i + b) \geq 1 - \xi_i, \xi_i \geq 0$$

$$\xi_i \geq 1 - y_i (w^T x_i + b) \quad \text{minimized when } \xi_i = 1 - y_i (w^T x_i + b)$$

$$\min_{w, b} \frac{1}{2} \|w\|^2 + \frac{c}{n} \sum (1 - y_i (w^T x_i + b))_+ = \min_{w, b} \frac{1}{2} \|w\|^2 + \frac{c}{n} \sum [1 - y_i (w^T x_i + b)]_+$$

l^2 -reg. penalty hinge loss