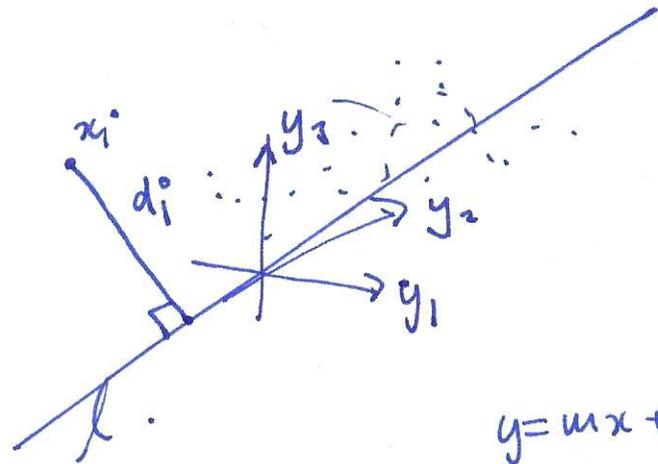


PCA: Principal component analysis

intuition

$$X = (x_1, \dots, x_n) \subseteq \mathbb{R}^d.$$

$$\uparrow$$
$$(x_{11}, x_{12}, x_{13}, \dots, x_{1d})$$



$$y = mx + c$$
$$a_0 + a_1 x + a_2 y = 0$$

- find the best fit line.

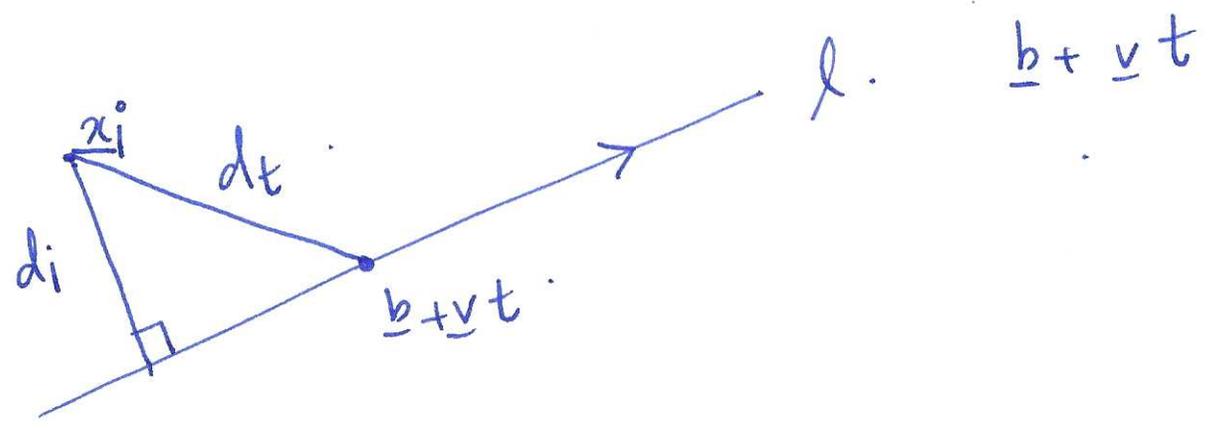
$$a_0 + a_1 y_1 + a_2 y_2 + \dots + a_d y_d = 0$$

choose a_i to minimize

sum of squares error.

$$\frac{1}{n} \sum_{i=1}^n d_i^2$$

can find d_i just from vector algebra:



$$\|d_t\|^2 = \|\underline{b} + \underline{v}t - \underline{x}_i\|^2 = (\underline{b} + \underline{v}t - \underline{x}_i) \cdot (\underline{b} + \underline{v}t - \underline{x}_i)$$

choose t to minimize this.

$$\frac{d}{dt} \|d_t\|^2 = 2 \underline{v} \cdot (\underline{b} + \underline{v}t - \underline{x}_i) = 0$$

$$\underline{v} \cdot \underline{b} + \underline{v} \cdot \underline{v}t - \underline{v} \cdot \underline{x}_i = 0$$

$$t = \frac{\underline{v} \cdot (\underline{x}_i - \underline{b})}{\|\underline{v}\|^2}$$

can do this for planes etc.

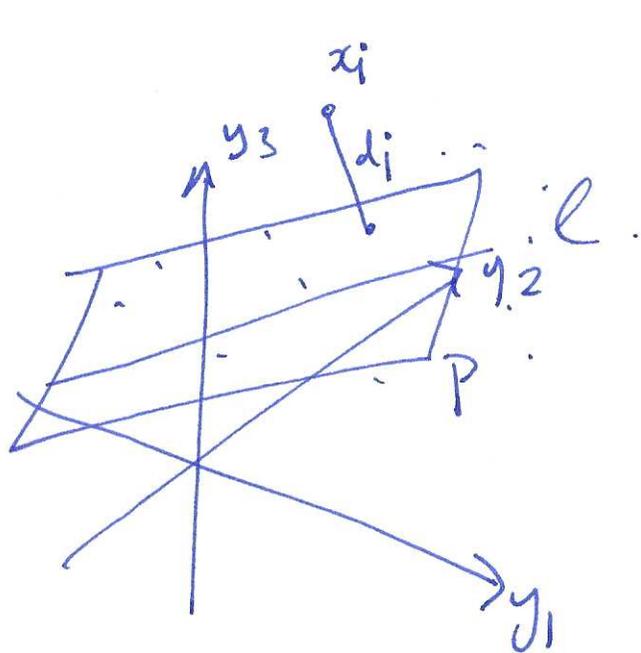
choose P to minimize

$$\frac{1}{n} \sum d_i^2$$

PCA (4 vectors)

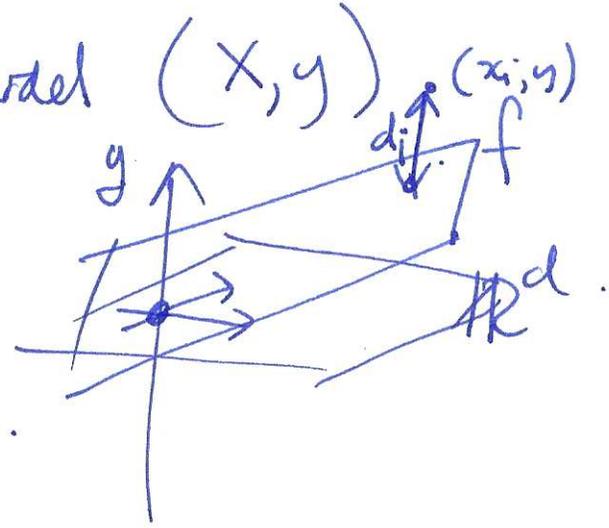
— vectors

— weights



Lasso recall linear regression model (X, y)

$X = (x_i) \in \mathbb{R}^d$



$f(x_i) = a_0 + a_1 x_{i1} + a_2 x_{i2} + \dots + a_d x_{id}$

choose a_0, a_1, \dots, a_d to minimize sum of

squares error $\frac{1}{n} \sum_{i=1}^n d_i$

some a_i 's will be small.

Lasso choose $\lambda = 0.01$
minimise

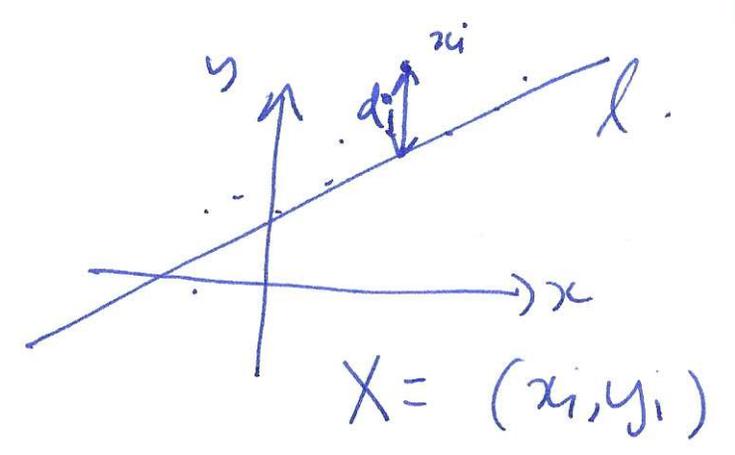
$\frac{1}{n} \sum_{i=1}^n d_i + \lambda |\# \text{ cols}|$

fix $\lambda = 0.01$

\leadsto get best a_0, a_1, \dots, a_k , rest are zero.
best k -dim linear model.

Lasso	vs	PCA
↑		↑
most important columns.		best fit directions

Bootstrap (type of cross-validation)



recall linear regression.

$$y = mx + b.$$

- get estimates for m, b .

- get accuracy / standard ^{error} estimates for m, b .

problem don't get error estimates for any

non-parametric model / many parametric models

solution : bootstrap.

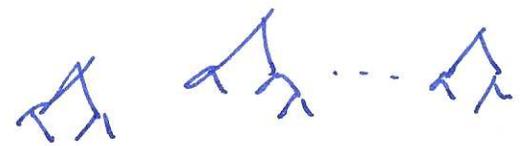
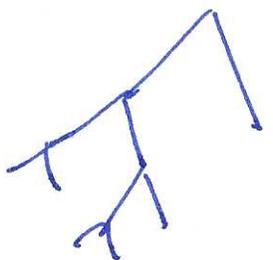
bootstrap • take lots of random samples from the original data set, compute estimates for these.

$$X = (x_1, \dots, x_n) \sim N_{\substack{15000 \\ 15000}}$$

$$B_1, B_2, \dots, B_{\substack{100 \\ 1000}} \leftarrow \text{of size } k \leq n$$

- take variance/standard deviation of these estimate
- in practice this is a good estimate of the variance of the model.

decision trees · bagging / bootstrap aggregation · boostings · random forest ⁽⁸⁾



- high variance.
- easy to understand.

bagging

idea : • subsample b times.

$$X \supseteq B_1, B_2, \dots, B_b.$$

- build a decision tree for each B_i
- take average / majority vote of the b trees.
(numeric) (categorical)
- tends to be lower variance than a single tree.

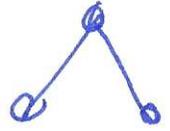
decision trees: boosting data set X . want $f(x_i) = y_i$ ⑨

parameters: b # trees.

$\lambda = 0.01$ (shrinkage parameter).

assume rescaled data.

d # splits. (# leaves = $d+1$) $(d=i)$.



① start with $f(x) = 0$ and predictions/residuals $r_0 = y$ (X, y).

② for $i=1$ to b

• fit a tree f^b with d splits to the data (X, r_{b-1}) .

• update. new $f_b := f_{b-1}(x) + \lambda f^b(x)$.

• update the residuals $r_i = r_{i-1} - \lambda f^b(x_i)$.

③ repeat ~~4~~ b times.