# Linear models

## Max-likelihood & least-squares

$$(*) \rightarrow y = w^T x + \epsilon, \quad \text{where} \quad \epsilon \sim N\left(0, \frac{1}{\beta}\right)$$

($w$: unknown)

($\beta$: precision parameter)

Suppose $\mathcal{D}_n = \{(x_i, y_i), i=1\text{--}n\}$ iid & satisfying $(*)$.

$$L(w, \beta) = \prod_{i=1}^{n} \frac{\sqrt{\beta}}{\sqrt{2\pi}} \exp\left(-\frac{\beta}{2}\left(y_i - w^T x_i\right)^2\right)$$

$$l(w,\beta) = \log L(w, \beta) = \frac{n}{2}\log\beta - \frac{n}{2}\log 2\pi - \beta \boxed{\frac{1}{2}\sum_{i=1}^{n}\left(y_i - w^T x_i\right)^2}$$

To maximize $l(w, \beta)$ wrt $w$, it is enough to

$$\hat{w}_{ML} = \arg\min_{w} \frac{1}{2}\sum_{i=1}^{n}\left(y_i - w^T x_i\right)^2 \quad \leftarrow \text{empirical risk minimization}$$

## Linear Regression:

Given data $\{(x_i, y_i); i=1\text{--}n\}$

$$x_i \in \mathbb{R}^d, \quad y \in \mathbb{R}$$

$$J(w) = \frac{1}{2}\sum_{i=1}^{n}\left(x_i^T w - y_i\right)^2$$

Minimize $J$:

$$A = \begin{bmatrix} - x_1^T - \\ - x_2^T - \\ \vdots \\ - x_n^T - \end{bmatrix} \quad n \times d$$

$$AW = \begin{bmatrix} x_1^T w \\ \vdots \\ x_n^T w \end{bmatrix} \qquad Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

$$AW - Y = \begin{bmatrix} x_1^T w - y_1 \\ \vdots \\ x_n^T w - y_n \end{bmatrix}$$

$$(AW-Y)^T (AW-Y) = \sum_{i=1}^{n} (x_i^T w - y_i)^2 = 2J(w)$$

$$\nabla_w J(w) = \frac{1}{2} \nabla_w (AW-Y)^T (AW-Y) = A^T (AW-Y)$$

So, $\quad \nabla J(w) = 0 \iff (A^T A) w = A^T Y$

Can we write $\quad w = (A^T A)^{-1} A^T Y$ ?

Yes, if $\quad$ A is full rank.

$\quad$ since $\quad$ full rank $A \Rightarrow A^T A$ is invertible

$\quad$ ( why? argue $rank(A) = rank(A^T A)$ by showing

$$N(A) = N(A^T A) \,)$$

Augmented features: $\{ (\tilde{x}_i, y_i), i=1 \dots n \} \quad \tilde{x}_i \in \mathbb{R}^d$

$$x_i = (1, \overleftarrow{\tilde{x}_i})$$

$$f(x) = \sum_{i=1}^{d} w_i x_i + w_o = \underset{\nwarrow}{w^T} \underset{\nearrow}{\tilde{x}} + w_o$$

$$d - \text{dimensional objects}$$

Why $w_0$?

$$J(w) = \frac{1}{2} \sum_{i=1}^{n} (w^T \tilde{x}_i + w_0 - y_i)^2$$

Take partial derivative wrt $w_0$

$$\frac{\partial J}{\partial w_0} = 0$$

$$\sum_{i=1}^{n} (w^T \tilde{x}_i + w_0 - y_i) = 0$$

Simplifying, $\quad w_0 = \frac{1}{n} \sum_{i=1}^{n} y_i - w^T \left( \sum_{i=1}^{n} \tilde{x}_i \right)$

---

Polynomial regression:

$$d=1, \quad \{ (x_i, y_i), i=1 \ldots n \} \quad x_i, y_i \in \mathbb{R}$$

use               transformed features, i.e.,

the following model

$$\hat{y}(x) = \underbrace{w_0 + w_1 x + w_2 x^2 + \cdots + w_m x^m}_{\text{mth degree polynomial}}$$

features $\leftarrow$ polynomial basis.

$$\hat{y}(x) = \sum_{j=0}^{m} w_j \phi_j(x), \quad \phi_j(x) = x^j$$

$$= w^T \phi(x)$$

$$w = (w_0, \text{---} \ w_m)$$

$$\phi(x) = (\phi_0(x), \text{----} \ \phi_m(x))$$

$$= (1, x, x^2, \text{---}, x^m)$$

Alternately, $\quad \phi_j(x) = \exp\left(-\dfrac{(x-\mu_j)^2}{2s^2}\right)$ ← Gaussian basis



$M=1$
$M=0$
$M=9$

Training error

Complexity (M)

Test Error

sweet spot

$M \rightarrow$

---

Least-squares: Geometric viewpoint

Recap of projections:-



$$p = \hat{x} \, a$$

$$(b - \hat{x} a) \perp a$$

$$a^T(b - \hat{x}a) = 0$$

$$\hat{x} = \frac{a^T b}{a^T a}$$

$$p = \hat{x}a = \left(\frac{a^T b}{a^T a}\right) a$$

$$= a\frac{a^T b}{a^T a} = \left[\frac{1}{a^T a}(aa^T)\right] b$$
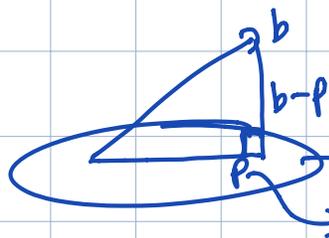
Projection matrix

$$P = \frac{1}{a^T a} aa^T$$

Projection matrix is (i) symmetric

(ii) $P^2 = P$

---

## Project onto a subspace?

A is a $m \times n$-matrix.

Want: Project b onto Col(A)



Col(A) = span (columns of A)

$$p = A\hat{x}$$

$$e = b - p = b - A\hat{x}$$

$e \perp$ every vector in C(A)

$\Rightarrow \quad c \in N(A^T) \quad$ because $\quad N(A^T) \perp C(A)$

$$A^T(b - A\hat{x}) = 0$$

$$A^T A \hat{x} = A^T b$$

observe $\quad A = \begin{bmatrix} & | & & & | \\ & a_1 & - & - & - & a_n \\ & | & & & | \end{bmatrix}$

$\left\{ \begin{array}{l} a_1^T(b - A\hat{x}) = 0 \\ \vdots \\ a_n^T(b - A\hat{x}) = 0 \\ \begin{bmatrix} a_1^T \\ \vdots \\ a_n^T \end{bmatrix} [b - A\hat{x}] = 0 \\ A^T(b - A\hat{x}) = 0 \end{array} \right.$

So, $\quad A^T A \hat{x} = A^T b \quad$ & this minimizes

$$E = \|Ax - b\|^2$$

If $\quad$ A is full col-rank, then

$$\hat{x} = (A^T A)^{-1} A^T b$$

Projection $\quad p = A\hat{x} = A(A^T A)^{-1} A^T b$

---

## Example:-

$$\begin{bmatrix} -1 & 1 \\ 1 & 1 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} \theta' \\ \theta'' \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 3 \end{bmatrix}$$

Want to solve $\quad A\theta = b$

Check if b is in C(A).

$$\begin{bmatrix} -1 & 1 & | & 1 \\ +1 & 1 & | & 1 \\ 2 & 1 & | & 3 \end{bmatrix} \rightarrow \cdots \cdots \cdots \rightarrow \begin{bmatrix} 1 & -1 & | & -1 \\ 0 & 1 & | & 1 \\ 0 & 0 & | & 2 \end{bmatrix}$$
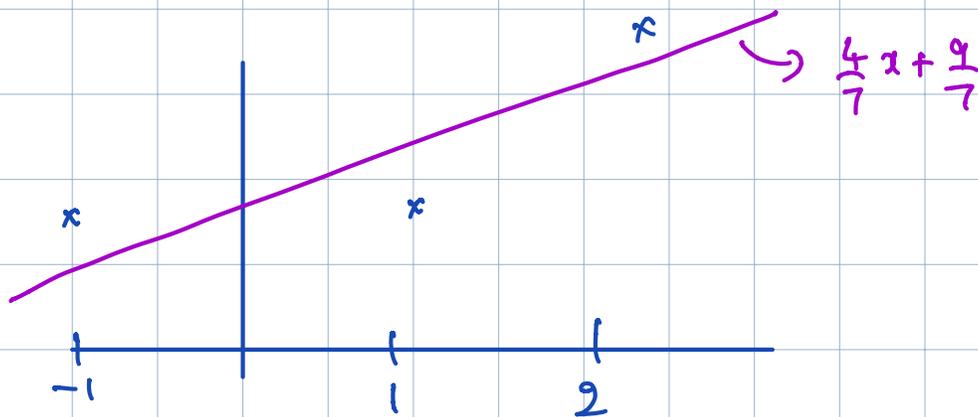
↗ inconsistent system

So $b \notin C(A)$.

Least-squares:  $\qquad A^T A \,\hat{\theta} = A^T b$

$$A^T A = \begin{bmatrix} 6 & 2 \\ 2 & 3 \end{bmatrix}$$

$$\begin{bmatrix} 6 & 2 \\ 2 & 3 \end{bmatrix} \begin{bmatrix} \hat{\theta}^1 \\ \hat{\theta}^0 \end{bmatrix} = \begin{bmatrix} 6 \\ 5 \end{bmatrix}$$

$$\hat{\theta}^1 = \frac{4}{7} \qquad\qquad \hat{\theta}^0 = \frac{9}{7}$$



$\rightarrow \frac{4}{7}x + \frac{9}{7}$

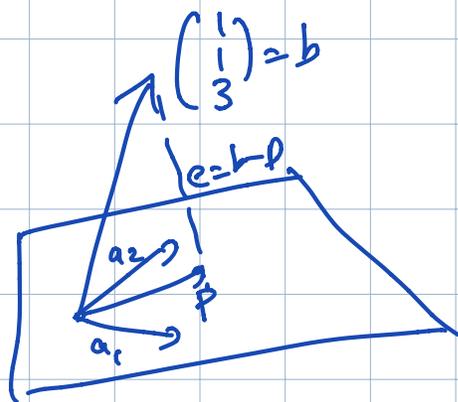$$P_1 = \frac{6}{7}, \quad P_2 = \frac{13}{7}, \quad P_3 = \frac{17}{7}$$

$$E^2 = \|b - A\hat{\theta}\|^2$$

$$= \left[1 - \left(-\frac{4}{7} + \frac{9}{7}\right)\right]^2 +$$

$$+ \left[1 - \left(\frac{4}{7} + \frac{9}{7}\right)\right]^2 + \left[3 - \left(\frac{8}{7} + \frac{9}{7}\right)\right]^2$$

$$= +\frac{2}{7}$$

$$e = b - p = \left(\frac{2}{7}, -\frac{6}{7}, \frac{4}{7}\right)$$



$$e \perp (-1, 1, 2)$$
$$a_1$$

$$e \perp (1, 1, 1)$$
$$a_2$$

---

## BIAS- VARIANCE  TRADEOFF

$f(\cdot) \to$ predictor   i.e.,  $f(x)$ is the prediction
for input feature $x$, & $y$ is the target
Suppose $(x, y)$ are chosen uzing some distribution $D$.
Then,  the  expected  loss  or  the "risk" is

$$R(f) = E\left((f(x) - y)^2\right)$$

$$\hookrightarrow \text{ over } x, y$$

**Claim:** $f^*(x) = E[y|x]$ is the best predictor,

ie., $f^*$ minimizes $R(f)$.

**Pf:** $E\left((f(x) - y)^2\right)$

$$= E\left[\left(f(x) - E[y|x] + E[y|x] - y\right)^2\right]$$

$(*)$ $\quad E\left(f(x) - y)^2\right) = E\left(f(x) - E[y|x]\right)^2 + E\left[\left(E[y|x] - y\right)^2\right]$

$\uparrow$

Predictor $f$ present only in this term
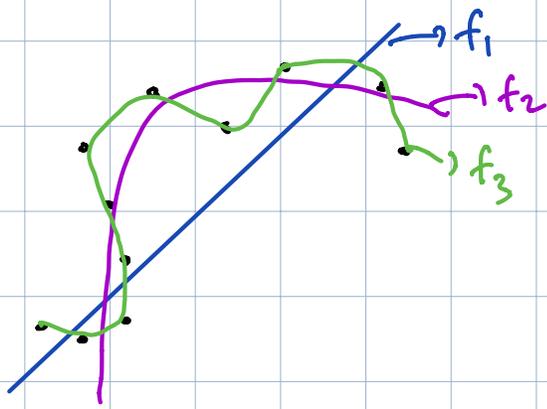
noise term

$(*)$ is minimized for $f = E[y|x]$ ▢

In a typical ML setting, $R(f)$ cannot be evaluated for a given $f$, since the underlying distributions are unknown.

So, collect training data $\{(x_i, y_i), i = 1 \cdots n\}$ sampled iid from $\mathcal{D}$, and minimize

Empirical $\rightarrow$ $R_n(f) = \dfrac{1}{n} \sum\limits_{i=1}^{n} (f(x_i) - y_i)^2$
risk

Intuitively,

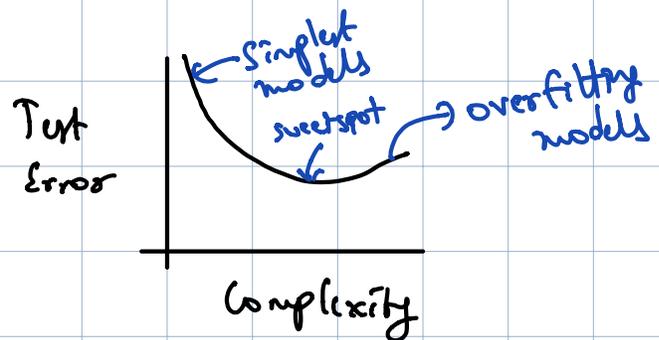$f_1$ is very simple

$f_3$ is very accurate

Maybe $f_2$ is the right fit

So, it is not enough to use $R_n(f)$ to judge $f$, since $f_3$ minimizes $R_n(f)$ (& fits noise). This phenomenon is referred to as "overfitting".

Test error = $\frac{1}{m} \sum_{i=1}^{m} \left( \bar{y}_i - f(\bar{x}_i) \right)^2$, where
(or generalization)
for $f$

$\{ (\bar{x}_i, \bar{y}_i), i=1 - m \}$ is the test data generated using distribution $\mathcal{D}$ (used for generating training data as well).



Training error

Complexity
(e.g., degree of the polynomial)



Test Error

← simplest models

sweetspot → overfitting models

Complexity

PTO

**Goal:** $\min_{f} E\left( \left( f_D(x) - y \right)^2 \right)$

$f_D(\cdot) \to$ predictor learnt using some dataset $D$.

$$E\left( \left( f_D(x) - y \right)^2 \right)$$

$$= E\left( \left( f_D(x) - E[y|x] \right)^2 \right) + E\left( \left( E(y|x) - y \right)^2 \right)$$

$\downarrow$
(I)

noise term
(unavoidable)

$$(I) = E\left( \left( f_D(x) - E(y|x) \right)^2 \right)$$

$$= E\left[ \left( f_D(x) - E_D(f_D(x)) \right)^2 \right.$$

$$+ \left( E_D(f_D(x)) - E(y|x) \right)^2$$

$$\left. + 2\left( f_D(x) - E_D(f_D(x)) \right)\left( E_D(f_D(x)) - E(y|x) \right) \right]$$

$E_D(f_D(x))$
$= E(f_D(x)|x)$
Fix $x$ &
average over
many datasets,
say $D_1, D_2 \text{---}$

$$= E\left[ \left( f_D(x) - E_D(f_D(x)) \right)^2 \right] \rightsquigarrow \text{Variance}$$

$$+ E\left[ \left[ E_D(f_D(x)) - E(y|x) \right]^2 \right] \rightsquigarrow (\text{Bias})^2$$

---

It is easy to see that as $f$ grows complex, the bias decreases.

**Claim:** As $f$ grows complex, variance increases.

"Curse of dimensionality".



100 points



$\leftarrow 1m^2$        $10^k$ points

\# points required to sample a unit hypercube
grows exponentially with the dimension

$$\min_{f \in \mathcal{F}_1} \hat{R}(f) \quad , \quad \min_{f \in \mathcal{F}_2} \hat{R}(f), \; \_\_ \; \& \; \text{so on, where}$$

$\mathcal{F}_i =$ set of all polynomials with
degree at most $i$.

Vector of co-efficients of a polynomial in
$\mathcal{F}_i$ sit in $\mathbb{R}^{d+1}$

With increasing $i$, one need to explore more points
to find the best $f$ in $\mathcal{F}_i$

(OR)

Given a fixed \# of points, the parameter
space is explored less efficiently for higher order $\mathcal{F}_i$,
leading to errors.

A work around: Add a penalty cost, i.e.,
solve the following variant of ERM:

$$\frac{1}{2} \sum_{j=1}^{n} (y_j - f(x_j))^2 + \lambda \underset{\uparrow}{L(i)} \quad —(\ast)$$

complexity cost

From ML discussion before,

$$\min \frac{1}{2} \sum_{j=1}^{n} (y_j - f(x_j))^2 \iff \max \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi}} \exp\left(\frac{(y_i - f(x_i))^2}{2}\right)$$

where $y_i = f(x_i) + \epsilon_i$

$\hookrightarrow$ standard Gaussian

Similarly, $(\ast)$ can be viewed as

$$\min \frac{1}{2} \sum_{j=1}^{n} (y_j - f(x_j))^2 + \lambda L(i)$$

$\Updownarrow$

$$\max \prod_{j=1}^{n} \frac{1}{\sqrt{2\pi}} \exp\left(- \frac{(y_j - f(x_j))^2}{2}\right) \times \underbrace{C(\lambda) \exp\left(-\lambda L(i)\right)}$$

$\downarrow$

prior probability on $\epsilon_i$

A choice for $L(i)$: $\quad L(i) = \| \beta(i) \|^2$

Regularized version of regression:- (Ridge regression)

$$\hat{R}_n(f) = \min_{w} \frac{1}{2} \sum_{i=1}^{n} (w^T x_i - y_i)^2 + \lambda \|w\|^2$$

Too small a $\lambda \rightarrow$ no effect of regularization
(overfit)

Too large a $\lambda \rightarrow$ under fit

Ref: Table 1.2

---

## Illustration of bias-variance tradeoff for a linear model:

Consider the model $y = w^T x + \epsilon$, $\epsilon \sim N(0, \frac{1}{\beta})$

The ML estimate $\hat{w}_{ML}$ for $w$, given $\mathcal{D} = \{(x_i, y_i), i=1,\dots,n\}$

$$\hat{w}_{ML} = (A^T A)^{-1} A^T Y, \text{ where}$$

$$A = \begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix}, \quad Y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}, \quad \eta = \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{pmatrix}$$

"Assume cols of $A$ are linearly independent".

(a) $\quad E(\hat{\omega}_{ML}) = E((A^TA)^{-1}A^TY)$

$\qquad\qquad\qquad = (A^TA)^{-1}A^T EY$

$\qquad\qquad\qquad = (A^TA)^{-1}A^T AW$

$\qquad\qquad\qquad = W.$

(b) $\quad Var(\hat{\omega}_{ML})$

$\quad = E((\hat{\omega}_{ML} - w)(\hat{\omega}_{ML} - w)^T)$

$\quad = E\left[((A^TA)^{-1}A^TY - w)((A^TA)^{-1}A^TY - w)^T\right]$

$\quad = E\left[((A^TA)^{-1}A^TY - w)(Y^TA(A^TA)^{-1} - w^T)\right]$

$\quad = (A^TA)^{-1}A^T E[YY^T] A(A^TA)^{-1} - ww^T$

$\quad = (A^TA)^{-1}A^T E\left[(AW+\eta)(AW+\eta)^T\right] A(A^TA)^{-1}$
$\qquad\qquad\qquad\qquad\qquad - ww^T$

$\quad = (A^TA)^{-1}A^T \left[ AWW^TA^T + \frac{1}{\beta}I_{n\times n}\right] A(A^TA)^{-1}$
$\qquad\qquad\qquad\qquad\qquad - ww^T$

$\quad = ww^T + \frac{1}{\beta}(A^TA)^{-1} - ww^T$

$\quad = \frac{1}{\beta}(A^TA)^{-1}$

(c) Bias-variance decomposition

$$E\left((f_D(x)-y)^2\right)$$

$$= E_{xy}\left(\underbrace{(E(y|x)-y)^2}_{noise}\right) + E_x\left[\left(\underbrace{E_D(f_D(x))-E(y|x)}_{bias}\right)^2\right]$$

$$+ E\left(\underbrace{f_D(x)-E_D(f_D(x))^2}_{Variance}\right)$$

for linear case, $\quad f_D(x) = \hat{w}_{ML}^T x$

$$Bias = \quad E_D(\hat{w}_{ML}^T x) - w^T x \quad = 0 \quad \left(\begin{array}{c}from \\ part(a)\end{array}\right)$$

$$Variance = \quad E\left((f_D(x) - E_D(f_D(x)))^2 \mid x\right)$$

Let
$\tilde{E} = \boxed{E(\cdot|x)}$

$$= \quad \tilde{E}\left((x^T\hat{w}_{ML} - x^T w)^2\right)$$

$$= \tilde{E}\left((x^T(A^TA)^{-1}A^T Y - x^T w)^2\right)$$

$$= \tilde{E}\left((x^T(A^TA)^{-1}A^T(Aw+\eta) - x^Tw)^2\right)$$

$$= \tilde{E}\left((x^Tw + x^T(A^TA)^{-1}A^T\eta - x^Tw)^2\right)$$

$$= \widetilde{E}\left( x^T (A^T A)^{-1} A^T \eta \right)^2 )$$

$$= \widetilde{E}\left( \left( x^T (A^T A)^{-1} A^T \eta \right) \left( x^T (A^T A)^{-1} A^T \eta \right)^T \right)$$

$$= x^T (A^T A)^{-1} A^T \underbrace{E\left( \eta \eta^T \right)}_{= \frac{1}{\beta} I_{n \times n}} \left( x^T (A^T A)^{-1} A^T \right)^T$$

$$= \frac{1}{\beta} x^T (A^T A)^{-1} A^T A (A^T A)^{-1} x$$

$$= \frac{1}{\beta} x^T (A^T A)^{-1} x$$

**H.w.**

① Let $C = A (A^T A)^{-1} A^T$

Show that $C$ is a projection matrix, i.e., $C$ is symmetric & $C^2 = C$

② Check if $(I - C)$ is a projection matrix

③ Let $S = \text{span} ( \text{cols of } A )$. Show that, for any $z \in \mathbb{R}^d$, $Cz$ is the projection of $z$ onto $S$.

④ Show that $A \hat{w}_{mL}$ is orthogonal to $Y - A \hat{w}_{mL}$.

---

**H.w.** Redo the exercise for a ridge regression-based

estimate $\hat{w}_{Reg}$.

(i) Write out the expression for $\hat{w}_{Reg}$

given $\mathcal{D} = \{ (x_i, y_i), i = 1-n \}$

(ii) $E\left( \hat{w}_{Reg} \right), \quad Var\left( \hat{w}_{Reg} \right)$

(iii) Calculate the bias & variance components.