

Linear Regression in Python

Objective: find line of best fit through a set of points

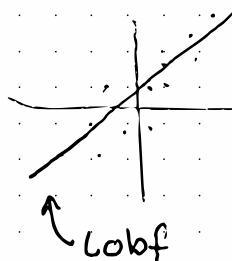
↳ eventually add more variables

What is ML: predict outcomes from past examples

- supervised ($X \rightarrow Y$)

- unsupervised (X) ↲ learn structure of the data

Supervised ↗ Classification
↘ Regression



Correlation ≠ causation

↳ randomised control trials

↳ synthesis of many studies (meta-studies)

example of linear equation. $\boxed{V = IR}$

\uparrow voltage
 \uparrow current
resistance

where measuring the V & I is what we are changing

↳ we can figure out R by taking measurements & finding slope

If we have \vec{x} & \vec{y} where const relating \vec{x} to \vec{y} is \hat{y}

$$\text{ie } \boxed{\hat{y}_i = a x_i + b}$$

How do we think of error between \vec{y} & \hat{y} ?

$$E^2 = \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

↳ No, error would be 0 when it's not on line

∴ take abs or sq error

$$MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

$$MAE = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i|$$

How do we optimise these equations?

$$E = \sum_{i=1}^N [y_i - (ax_i + b)]^2$$

needs to
find deriv.
wrt a & b

$$\frac{\partial E}{\partial a} = \sum_{i=1}^N 2(y_i - (ax_i + b)) * (-x_i)$$

optimise, so $\frac{\partial E}{\partial a} = 0$

$$\sum_{i=1}^N -y_i x_i + ax_i^2 + bx_i = 0$$

$$\sum_{i=1}^N ax_i^2 + bx_i = \sum_{i=1}^N y_i x_i$$

$$\frac{\partial E}{\partial b} = 0 = \sum_{i=1}^N 2(y_i - (ax_i + b)) * (-1)$$

$$\sum_{i=1}^N ax_i + bN = \sum_{i=1}^N y_i$$

doesn't seem to work...

Solve using Linear Algebra

or

$$a = \frac{\bar{xy} - \bar{x}\bar{y}}{\bar{x}^2 - \bar{x}^2}$$

$$b = \frac{\bar{y}\bar{x}^2 - \bar{x}\bar{xy}}{\bar{x}^2 - \bar{x}^2}$$

Notation is not clear

Note that $\bar{a} \cdot \bar{b} = \vec{a}^T \vec{b}$

How did we arrive at these results:

1. define equation to estimate
2. define cost; defines "hyperplane"
3. find parameters that minimise cost.
 \hookrightarrow lowest cost = solution

Note: optimization cannot work here
this problem is not convex

\rightarrow Gradient descent

init w/ random vars a, b

\hookrightarrow predict $\hat{y} = ax + b$

find loss = $\frac{1}{N} \sum_{i=1}^N (y - \hat{y})^2$

need to update based on error

$$\frac{\partial \hat{y}}{\partial a} \left[\frac{1}{N} \sum_{i=1}^N (y - ax - b)^2 \right]$$

$$= \frac{2}{N} \sum_{i=1}^N (y - ax - b) (-x)$$

$$\frac{\partial \hat{y}}{\partial b} = \frac{2}{N} \sum_{i=1}^N (y - ax - b) (-1)$$

Update based on grad.

$$a = \alpha * \frac{\partial \hat{y}}{\partial a}$$

$$b = \alpha * \frac{\partial \hat{y}}{\partial b}$$

yay! I can believe we did it!

R² values → defined as

$$R^2 = 1 - \frac{\sum_{i=1}^N (y - \hat{y})^2}{\sum_{i=1}^N (y - \bar{y})^2}$$

- $R^2 = 0$: prediction is the same as taking the mean
- $R^2 = 1$: prediction is perfect
- $R^2 = -1$: prediction gives more errors than taking the average

Multiple Linear Regression

↳ normally multiple effects influence the dependent variable

$$\hookrightarrow 1-d: \hat{y} = \Theta^T \vec{x} + b$$

$$\hookrightarrow \Theta_1 x_1 + \Theta_2 x_2 + \dots + \Theta_n x_n$$

↳ New optimal solution

$$\frac{\partial E}{\partial \theta_j} = \sum_{i=1}^N 2(y_i - \theta^T x_i)(-x_{ij}) = 0$$

↙ vector ↘ scalar

$$\boxed{Nb} \frac{-\partial(\theta^T x_i)}{\partial \theta_j}$$

All other θ are 0

because they are const. $\Rightarrow = -\underline{\partial(\theta_j x_{ij})}$

↳ expand: $\sum_{i=1}^N y_i(x_{ij}) - \sum_{i=1}^N \theta^T x_i(-x_{ij}) = 0$

We do this
to try & solve
for θ

$$\therefore \sum_{i=1}^N y_i x_{ij} = \sum_{i=1}^N \theta^T x_i x_{ij}$$

↑ independent
of i

equation exists for every θ_j

- | | | |
|---|--|--|
| 1 | $\theta^T \sum_{i=1}^N x_i x_{i1} = \sum_{i=1}^N y_i x_{i1}$ | |
| 2 | $\theta^T \sum_{i=1}^N x_i x_{i2} = \sum_{i=1}^N y_i x_{i2}$ | $\Rightarrow \boxed{\theta^T (X^T X) = Y^T X}$ |
| ⋮ | ⋮ | ⋮ |
| D | $\theta^T \sum_{i=1}^N x_i x_{iD} = \sum_{i=1}^N y_i x_{iD}$ | Hallelujah! |

$$\Theta^T(X^T X) = Y^T X \quad \checkmark$$

$$(X^T X) \Theta = X^T Y$$

$$\Theta = (X^T X)^{-1} X^T Y \quad \checkmark$$

How to solve multiple LR using matrices

Same as last section only entirely in matrix form.

$$S = \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

$$= (Y - \hat{Y})^T (Y - \hat{Y}) \quad Y = X\Theta$$

$$= Y^T Y - Y^T \hat{Y} - \hat{Y}^T Y + \hat{Y}^T \hat{Y}$$

$$= (X\Theta)^T (X\Theta) - (X\Theta)^T \hat{Y} - \hat{Y}^T (X\Theta) + (\hat{Y})^T \hat{Y}$$

$$\frac{\partial S}{\partial \Theta} = 0 = -2X^T \hat{Y} + 2X^T X \Theta$$

$$\therefore \Theta = (X^T X)^{-1} X^T \hat{Y}$$

Shape of Θ

$$\Theta = (X^T X)^{-1} X^T Y$$

Say $X \quad Y$
 $100 \times 2 \quad 100 \times 1$

$$I = X^T X = 2 \times 2$$

$$I^T = (2 \times 2, 2 \times 100) = 2 \times 100$$

$$(X^T X)^{-1} X^T Y = (2 \times 100, 100 \times 1) = \boxed{2 \times 1}$$

$$\therefore \Theta = 0 \times 1 \leftarrow \text{vector length of dimensions}$$

... What about constant?

↳ add extra column to X made of 1

Polynomial Regression

Why is this still linear?

What we predict is a linear combo of weights (Θ). Thus, x can change how we weight.

Not creating code for this but multiple APIs are capable of this calculation

Interpreting the weights

↳ $y = mx + b \leftarrow$ fudge factor that is always present.
 ↑ "for every increase in x ,
 y is incremented by m "

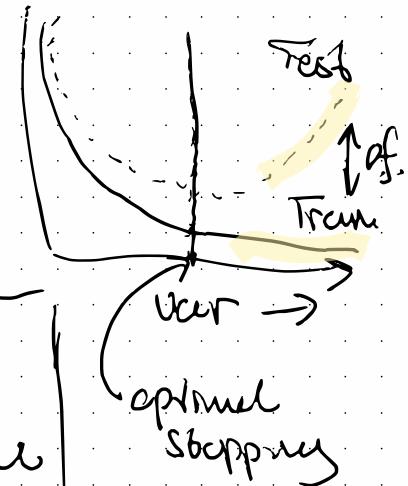
Generalization

Why not just fit data as closely as possible?

- ↳ doesn't generalize & noise is captured
- ↳ can't predict the future, only the past

Train on 80% of data

Test on 20% for
generalization



1. Train] True model
2. Test
3. User ← use this when model
is best fit yet

What about categorical vars?

i.e. degree type: B, M, P

IS-B	IS-M	IS-P
0	0	1
0	1	0
1	0	0

↓
bool values

Probabilistic Interpretation of sg. error

Show LR is max. likelihood solution to b.o.b

What is max. likelihood

() "What answer is most likely to be true"

$$\hookrightarrow E = \sum_{i=1}^n (y_i \cdot g_i)$$

ie, what is mean of Gaussian dist?

$$P(x_i) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right)$$

If we measure heights, we know they're independent & identically distributed

↳ not related & come from same dist.

$$\begin{aligned} \therefore P(X) &= \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(X_i - \mu)^2}{2\sigma^2}\right) \\ &= P(X|\mu) \end{aligned}$$

↑ prob of data given μ

We want to find μ such that $P(X)$ is max.

$$\frac{\partial P(X)}{\partial \mu} = 0? \quad \text{super hard as } N \uparrow$$

$$\therefore \frac{\partial \log(P(X))}{\partial \mu} = 0 \quad \checkmark$$

$$\text{so.. } L(X|\mu) = \log \left[\prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{-(x_i - \mu)^2}{2\sigma^2}\right) \right]$$

$$= \sum_{i=1}^N \left[-\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2} \frac{(x_i - \mu)^2}{\sigma^2} \right]$$

where $L(x) = \log[P(x)]$

$$\begin{aligned} \frac{\partial L(x|\mu)}{\partial \mu} &= \sum_{i=1}^N \frac{-2}{2} * \frac{(x_i - \mu)}{\sigma^2} * (-1) \\ &= \sum_{i=1}^N \frac{(x_i - \mu)}{\sigma^2} = 0 \end{aligned}$$

$$\text{solve for } \mu \dots \sum_{i=1}^N x_i = \sum_{i=1}^N \mu$$

$$\therefore \mu = \frac{1}{N} \sum_{i=1}^N x_i$$

we do this because we summed μ N times.

Moral of story?

w/ LR minimizing cost is equal to max.

log likelihood

(y is gaussian dist across independent vars.)

L2 regularisation "Ridge Regression"

MSE is vulnerable to outliers due to exp. in loss function. How to fix?

↳ penalize large weights

$$J = \sum_{i=1}^N (y_i - \hat{y}_i)^2 + \lambda |\boldsymbol{\theta}|^2$$

$= \boldsymbol{\theta}^T \boldsymbol{\theta}$

→ Probabilistically

$$P(Y|X, \boldsymbol{\theta}) = \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{1}{2\sigma^2} (y_n - \boldsymbol{\theta}^T \mathbf{x}_n)^2\right]$$

Both vars ↑
have norm. dist.

Now, we mix in a prior:

$$P(\boldsymbol{\theta}) = \frac{1}{K\pi} \exp\left[-\frac{1}{2} \boldsymbol{\theta}^T \boldsymbol{\theta}\right]$$

$$P(w|Y, X) \propto P(Y|X, w) * P(w)$$

↑ max a posteriori

So, how the fresh heck do we optimize this?¹⁷

(\Rightarrow find cost, take deriv & set to 0)

$$J = (Y - X\theta)^T (Y - X\theta) + \lambda \theta^T \theta$$

Note how we expect a similar form to before

$$\frac{\partial J}{\partial \theta} = -2X^T Y + 2X^T X\theta + 2\lambda\theta = 0$$

$$X^T X\theta + \lambda\theta = X^T Y$$

$$(X^T X + \lambda I)\theta = X^T Y$$

$$\boxed{\theta = (X^T X + \lambda I)^{-1} X^T Y}$$

I did A'!!

Dummy variable trap

2 types of encoding $\begin{cases} \rightarrow \text{One hot} \\ \rightarrow k-1 \end{cases}$ $\begin{matrix} \text{k classes} \\ = k \text{ vars} \\ \text{k classes} \\ = k-1 \text{ vars} \end{matrix}$

$k-1$ encoding puts one class into binary variable

Why? $X^T X$ becomes singular

\hookrightarrow columns no longer linearly indep.

→ How to deal with the trap?

- (\hookrightarrow) k-1 encoding, drop bias term, use $\|\cdot\|_2$, use gradient desc.

↑ used through deep learning.

Why doesn't L1 reg work?

- (\hookrightarrow) L1 can still be non-conv. indep.

→ Other problems: Huber loss

- (\hookrightarrow) 2+ vars are correlated to each other

Gradient descent

- (\odot) ← random setting for weights



$$\Theta := \Theta - \eta \nabla_{\Theta} J(\Theta)$$

↑ gradient vector of
 $J(\Theta)$

↑ take steps
down a mountain

$$\Theta := \Theta - (\eta * X^T (\hat{Y} - Y))$$

↑ updated weights

What is so important about L.D?

- ↳ used when loss function is not convex.
i.e. no global min that's easy to calculate

→ How to choose η ?

- ↳ too large: never converge
- ↳ too low: slower computation

L1 Regularisation "Lasso"

- ↳ used for feature selection

$$J = \sum_{n=1}^N (y_n - \hat{y}_n)^2 + \lambda \|\theta\|_1$$

- ↳ probabilistically y_i : $p(\theta) = \frac{1}{2} \exp(-\lambda |\theta|)$

↳ Laplace dist.

$$\frac{\partial J}{\partial \theta} = -2X^T y + 2X^T X \theta + \lambda \text{sign}(\theta) = 0$$

Where $\text{Sign}(x) \begin{cases} = 1 & \text{when } x > 0 \\ = 0 & \text{when } x = 0 \\ = -1 & \text{when } x < 0 \end{cases}$

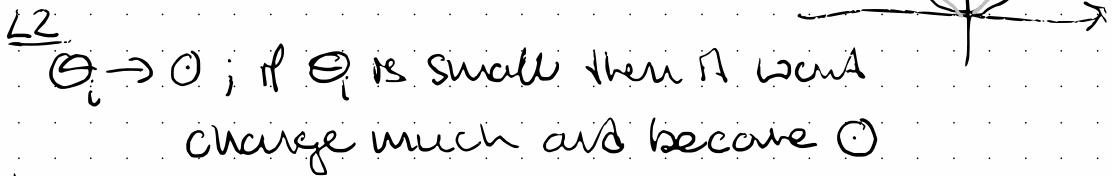
↑ can't
solve
for θ

L1 vs L2 reg

→ Help to prevent overfitting

- ↳ L1 chooses most important features
- ↳ L2 doesn't allow for huge weights.

This is because L2 penalty is quas & L1 is an abs func



L2

$\Theta_i \rightarrow 0$; if Θ_i is small then it won't change much and become 0

L1

If $\Theta_i = 0$ then it will stay there forever
(sparsity)

Combine L1 & L2 to create elastic net

$$J_{\text{en.}} = J + \lambda_1 |w| + \lambda_2 |w|^2$$