# Classifier evaluation

- Test set carved out of training data

# Evaluation metric?

Accuracy — percentage of correct answers

Problem — unbalanced categories

Often, the interesting case is a minority

- Fraud, Junk Mail, Rare disease

Suppose "Yes" occurs 5% of time

Blind "No" classifier is 95% accurate

Want to force classifier to flag "Yes"

Categorize errors more finely

|  | Prediction | |
|---|---|---|
| | Y | N |
| **Actual answer** Y | ✓ | ✗ |
| N | ✗ | ✓ |

5% Yes, Trivial "No" classifier

Y    N    goal

|       | Y | N |
|-------|---|---|
| Y     | 0 | 50 |
| N     | 0 | 950 |

will cause this

1000 cases
950 N    50 Y

$$\frac{TP}{TP+FP}$$

PRECISION

|   | Y | N |
|---|---|---|
| Y | True Positive | FN |
| N | FP | True Negative |

$$\frac{TP}{TP+FN}$$

Actually found

should have found

RECALL

# Precision - recall tradeoff

Screening test   vs   interview
Corona virus   vs   pancreatic cancer

Single number?

F-score : Harmonic mean

Reciprocal of mean of reciprocals

$$\frac{1}{\left(\dfrac{\frac{1}{P} + \frac{1}{R}}{2}\right)} \qquad \frac{2PR}{P+R}$$

Regression — predicting a numeric value

Fit a function to the data
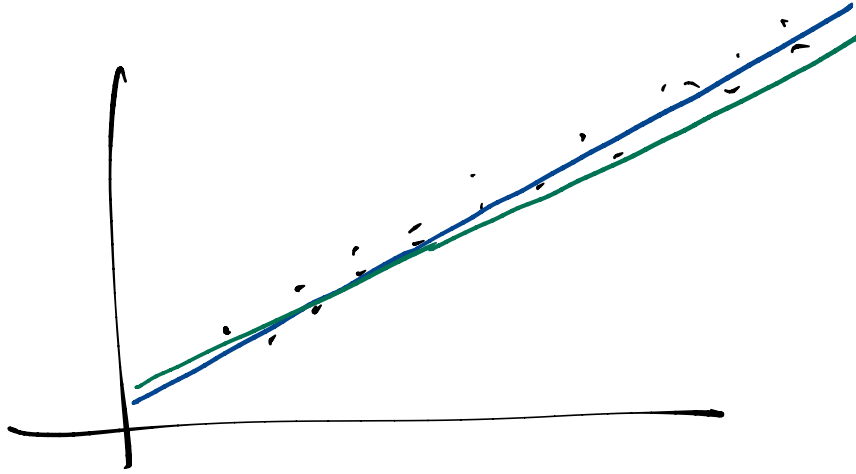
Attributes $A_1, \ldots, A_k$

$$f(x_1, \ldots, x_k)$$

Simplest case: $f(x_1 \ldots x_k) = a_1 x_1 + a_2 x_2 + \ldots a_k x_k + b$

$$f(x) = mx + b$$

How to find "best" $m, b$
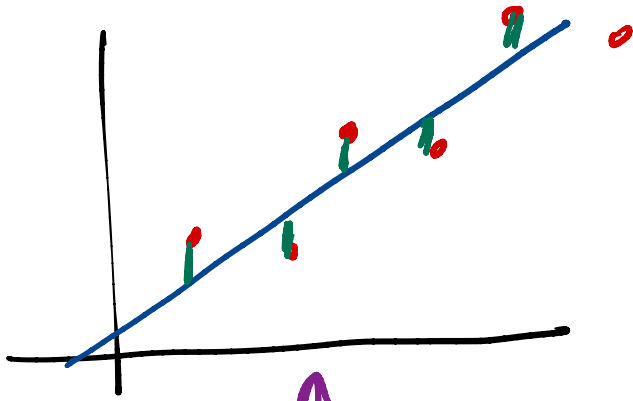
# Define an error measure

$(x_1, y_1)$

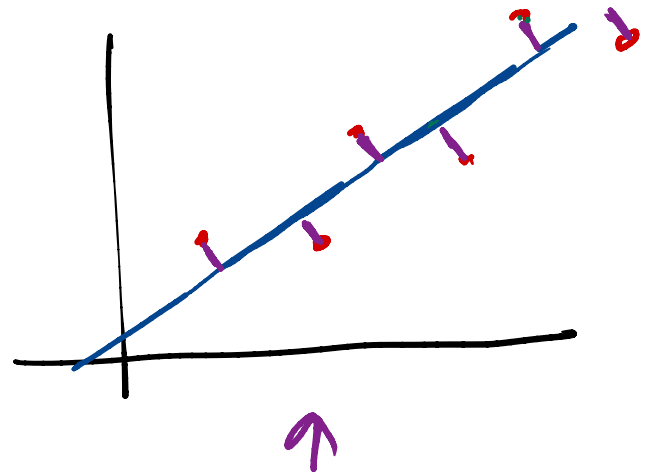$(x_2, y_2)$

$\vdots$

$(x_n, y_n)$

Prediction: $mx_i + b$

$$\left| (mx_i + b) - y_i \right|^2 \quad - \text{square error}$$

VS

↑
Square error

↑
Not this

$$\text{Mean Square Error} = \frac{1}{n} \sum_{i=1}^{n} \left( (mx_i + b) - y_i \right)^2$$

Find m, b to minimize MSE

Statistics — direct formula for $m, b$ based on
mean, variance etc of training points

Instead — iteratively improve $m$ & $b$

Adjust $m$ & $b$ so that MSE reduces

$\boxed{MSE}$ — Error, Loss, Cost $\quad \Theta(m, b)$

$$\Theta(m, b) = \frac{1}{n} \sum_{i=1}^{n} \left( (mx_i + b) - y_i \right)^2$$

## Want

$$\frac{\partial \Theta}{\partial m} \, , \, \frac{\partial \Theta}{\partial b}$$

Remember that $x_i$'s are fixed values, & $y_i$

$$2\frac{1}{n} \sum_{i=1}^{n} (\quad -\quad)^2$$

$$\frac{\partial \Theta}{\partial m} \quad 2\left(mx_i+b-y_i\right)\cdot x_i$$

$$\frac{\partial \Theta}{\partial b} \quad 2\left(mx_i+b-y_i\right)\cdot 1$$

Adjust m by $\quad \alpha \cdot \dfrac{-\partial \theta}{\partial m}$

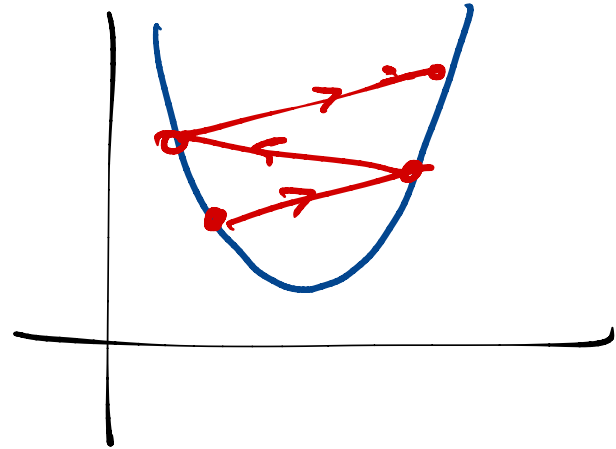$\qquad$ b by $\quad \alpha \cdot \dfrac{-\partial \theta}{\partial b}$

$\alpha$ is a small value

"learning rate"

If $\alpha$ is too small — progress is slow

If $\alpha$ is too big?

**Gradient descent** — batch of "predictions"
update coefficient

Can also do smaller batches & update incrementally

**Stochastic Gradient Descent (SGD)**

Pick a random subset to update

Recompute gradient

Repeat

Suppose the function is not linear?

Classically — transform the data

Suppose $f(x_1) = a_1 x_1 + a_2 x_1^2 + a_3$

$$a_1 x_1 + a_2 x_2 + a_3$$
$$\downarrow$$
$$x_2 = x_1^2$$

$$x_1, y_1 \longrightarrow x_1, x_1^2, y_1$$
$$x_2, y_2 \qquad\quad x_2, x_2^2, y_2$$
$$\vdots \qquad\qquad\quad \vdots$$
$$x_n, y_n \qquad\quad x_n, x_n^2, y_n$$
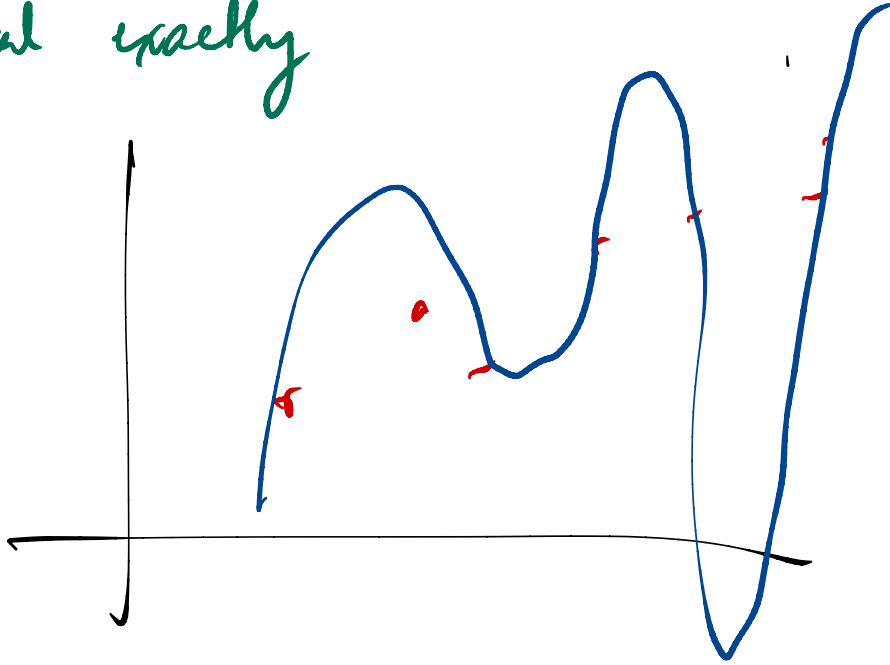
Python - sklearn

Specify degree of regression

Input is $(x_1, x_2, y)$

Degree is 2

$$(x_1, x_2, x_1^2, x_2^2, x_1 x_2, y)$$

We can always fit an arbitrarily high degree
polynomial exactly



Overfitting