New idea for a learning algorithm...

$$g = A(\mathbb{D}, \mathcal{H}), \quad \text{Predict via } \hat{y}^* = g(\vec{x}^*)$$

What if $g$ found the "closest" $\vec{x}_i \in \mathbb{D}$ to $\vec{x}^*$

and returned $\hat{y}^* = y_i$?   This closest $\vec{x}_i$ is called its neighbor.

Need to define "closest" via a distance function

$$d(\vec{x}_i, \vec{x}_k) = \| \vec{x}_i - \vec{x}_k \|_2^2 = (\vec{x}_i - \vec{x}_k)^\top (\vec{x}_i - \vec{x}_k) = \sum_{j=1}^{p} (x_{ij} - x_{kj})^2$$

Euclidean Norm Sqd.

many others...

$$\mathcal{H} = \{ \; ? \; \} \quad \Bigg\}  \text{ Difficult to define precisely}$$
$$A = ?$$

$g(\vec{x}^*)$ ... all the work happens here ...

Extension: find the $K$ closest $\vec{x}_i$'s. Then $\hat{y} = \text{Mode}[y_{(1)} ... y_{(K)}]$

↑ returns most likely class

$K$ nearest neighbors or "kNN" algorithm.

weaknesses?  $p$ large AKA "curse of dimensionality"

and not all $x_j$ terms are equally predictive!

choices?  $k$, $d$ — they really matter!! Learning is not simple.

So far, we have been concerned with problems s.t. $\mathcal{Y} = \{0,1\}$.
This is called "binary classification". If $\mathcal{Y} = \{0, 1, ..., K\}$
where the response levels are nominal (ie. no order),
this is called "classification" or "multilevel classification".

What if $\mathcal{Y} \in \mathbb{R}$ or $\mathcal{Y} \in R \subset \mathbb{R}$? This is then called "regression".

Why? For a historical reason which we will get to.

Can the threshold, perceptron or SVM do regression? Not without serious adaptation. What do we do?? Null Model? $\mathcal{H} = \{y_0 : y_0 \in \mathbb{R}\}$
$g = \bar{y}$ Avg. value

Recall $\mathcal{H} = \left\{ \mathbb{1}_{\vec{w} \cdot \vec{x}} : \vec{w} \in \mathbb{R}^{p+1} \right\}$

linear model

Why the indicator function? To coerce the output $g \in \{0,1\}$ instead of $\mathbb{R}$.
Why not use just the linear model?

conversion to index this at 0 and not call it b.

$\mathcal{H} = \{\vec{w} \cdot \vec{x} : \vec{w} \in \mathbb{R}^{p+1}\} = \{w_0 + w_1 x_1 + ... + w_p x_p : w_0 \in \mathbb{R}, w_1 \in \mathbb{R}, ..., w_p \in \mathbb{R}\}$

This is the most famous model period. The "linear regression model". For historical
reasons, we will denote the weights as $\vec{b}$, not $\vec{w}$.

$\mathcal{H} = \{\vec{b} \cdot \vec{x} : \vec{b} \in \mathbb{R}^{p+1}\} = \{b_0 + b_1 x_1 + b_2 x_2 + ... + b_p x_p : b_0 \in \mathbb{R}, b_1 \in \mathbb{R}, ..., b_p \in \mathbb{R}\}$

conversion

Dim of param space: $p+1$

We can visualize this if $p = 1$

$\vec{w}$ is the linear coefficients

$\mathcal{H} = \{w_0 + w_1 x_1 : w_0 \in \mathbb{R}, w_1 \in \mathbb{R}\}$

$$h^*(\vec{x}) = w_0^* + w_1^* x = \beta_0 + \beta_1 x$$

the candidate in $\mathcal{H}$ that most closely resembles $f$.

special notation for the optimal/true coefficients

Review: why doesn't $h^* = y$ exactly?

$$y = h^*(\vec{x}) + \varepsilon = h^*(\vec{x}) + \underbrace{(t(\vec{z}) - f(\vec{x}))}_{\text{ignorance}} + \underbrace{(f(\vec{x}) - h^*(\vec{x}))}_{\text{misspecification of linear model}}$$

$$\underbrace{\hspace{6cm}}_{\varepsilon}$$

AKA "noise" or "errors"

$h^*$ is inaccessible since we have to make an imperfect fit with finite data

$$y = g(\vec{x}) + e = f(\vec{x}) + \underbrace{(t(\vec{z}) - f(\vec{x}))}_{} + \underbrace{(f(\vec{x}) - h^*(\vec{x}))}_{\varepsilon} + \overbrace{(h^*(\vec{x}) - g(\vec{x}))}^{\text{estimation error}}$$

"residuals"  $\underbrace{\hspace{8cm}}_{e - \varepsilon}$

As $n \to \infty$  $g(\vec{x}) \to h^*(\vec{x})$  and $e - \varepsilon \to 0$  but  $y \neq g(\vec{x})$ since the other two errors are still present.

_____

AKA "simple regression"

Back to the linear model for $p=1$. How to fit $\vec{w}$, the parameters? Need an A.

First need a **loss** function, error function, objective function, cost function.

Recall before  $SSE := \overset{\min!}{\underset{i=1}{\sum^n}} \underbrace{(y_i - \hat{y}_i)^2}_{e_i} = \sum e_i^2$  "sum of squared error"

$$= \sum_{i=1}^n \left(y_i - \underbrace{(w_0 + w_1 x_{1i})}_{\bar{y} = \frac{1}{n}\sum y_i}\right)^2 = \sum_{i=1}^n y_i^2 + w_0^2 + w_1^2 x_i^2 - 2y_i w_0 - 2y_i w_1 x_i + 2 w_0 w_1 x_i$$

$$= \sum y_i^2 + n w_0^2 + w_1^2 \sum x_i^2 - 2n\bar{y} w_0 - 2 w_1 \sum x_i y_i + 2 w_0 w_1 n \bar{x}$$

Choose $w_0, w_1$ to min. the above

$$\frac{\partial}{\partial w_0}[\ ] = 2n w_0 - 2n\bar{y} + 2 w_1 n \bar{x} \overset{\text{set}}{=} 0 \implies \hat{w}_0 = \bar{y} - w_1 \bar{x}$$