L2

Recap :

①     representation

② search :

    Document $\Rightarrow$ vector of term frequency

      $X_i(j) = tf_{i(j)}$   "Term Frequency"

      $X_i(j) = tf_i(j) \cdot idf(j)$        "Inverse Document Frequency"

      $idf(j) = \log\left(\frac{n}{n_j}\right)$ $\longrightarrow$ number of documents

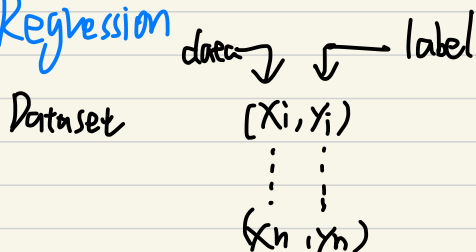                   $\searrow$ number of documents containing word $n_j$

| Pros | | Cons |
|------|---|------|
| — simple | | — $O(nd)$ time |
| | | per test instance |
| — robust | | |

Regression   data $\searrow$ $\swarrow$ label

Dataset       $(X_i, Y_i)$

           $\vdots$   $\vdots$

           $(X_n, Y_n)$

goal of regression : Find a function of   $Y_i \approx f(X_i)$
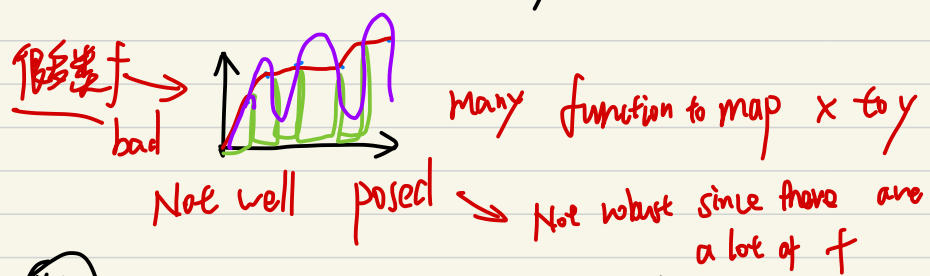
**Application :** ① Classification of Image :

$x_i$ image $y_i \rightarrow$ cat / no cat (类别)

② Retail Pricing : $x$: time  $y$: price of product

③ weather :

$x$: location $y$: rainfall

很复杂 → bad  Many function to map $x$ to $y$

Not well posed → Not robust since there are a lot of $f$

寻找好的线性模型 (Fix) Regression: find <u>a</u> function $f$ that belongs to a

class of functions $H$  固定一类映射函数

Hypothesis class ←

---

Linear models:

$H$ : set of linear functions

change in input ∝ change in output

**Why linear models?**
- simplicity
- stable behaviour
- easy to compute
- Interpretable
  高效、清晰
- building block of more complex function

$f(x) = f(x_0) + f'(x_0)(x-x_0) + \cdots$

泰勒展开, 线性模拟.

linear regression (univariate)

$x \rightarrow$ scalar $\qquad (x_1, y_1)$

$y \rightarrow$ scalar $\qquad (x_2, y_2)$

$\qquad\qquad\qquad\quad \vdots$

$\qquad\qquad\qquad (x_n, y_n)$

**Step 1: Representation**

$$y = w_0 + w_1 x$$

$$\hat{y_i} = w_0 + w_1 x_i$$

**Step 2: measure of goodness:**

eg. MSE ( mean square error )

$$E_x(\hat{y} - y)^2$$

测量值

predicted value

$$MSE = \frac{1}{n} \sum_{i=1}^{n} (w_0 + w_1 x_i - y_i)^2$$

寻找 $w_0$, $w_1$
使误差平方最小
关系

**Step 3: find best $w_0$, $w_1$ that minimize MSE**

take partial derivative over $w_0$, $w_1$    偏导

① $\dfrac{\partial MSE}{\partial w_0} = 0$

找最小MSE对应的
线参 $w_0$, $w_1$

$2 \cdot \dfrac{1}{n} \sum_{i=1}^{n} (w_0 + w_1 x_i - y_i) = 0$

solve for $w_0$

$$w_0 = \frac{\sum_i y_i}{n} - w_1 \frac{\sum x_i}{n}$$

$$w_0 = \bar{y} - w_1 \bar{x}$$

② $\dfrac{\partial MSE}{\partial w_1} = 0$

$2 \cdot \dfrac{1}{n} \sum_{i=1}^{n} x_i (w_0 + w_1 x_i - y_i) = 0$

$\dfrac{1}{n} \cdot \sum x_i (\bar{y} - w_1 \bar{x} + w_1 x_i - y_i) = 0$

$w_1 \cdot \dfrac{1}{n} \sum_{i=1}^{n} (x_i - \bar{x}) x_i = \dfrac{1}{n} \sum (y_i - \bar{y}) x_i$

协方差
$w_1 = \dfrac{\frac{1}{n} \sum (x_i - \bar{x})(y_i - \bar{y})}{\frac{1}{n} \sum (x_i - \bar{x})^2}$  方差

$Cov(X, Y)$
$= E[(x - \bar{x})(y - \bar{y})]$
$= E(xy) - E(x)E(y)$

$w_1 = \dfrac{(\frac{1}{n} \sum x_i y_i) - (\bar{x}\bar{y})}{\frac{1}{n} \sum x_i^2 - \bar{x}^2}$  variance of X

cross-covariance of X & y

**Variance**

$$\frac{1}{n} \sum_{i=1}^{n} (x_i - \bar{x})^2$$

$$= \frac{1}{n} \sum (x_i^2 + \bar{x}^2 - 2x_i \bar{x})$$

$$= \boxed{\frac{1}{n} \sum x_i^2 - \bar{x}^2}$$

$\sigma_x^2 \rightarrow$ variance of $x$

$\sigma_y^2 \rightarrow$ variance of $y$

$\sigma_{xy} \rightarrow$ cross-covariance of $x$ & $y$

$$MSE = \frac{1}{n} \sum_{i=1}^{n} \left[ y_i - (\underbrace{\bar{y} - w_1 \bar{x}}_{w_0} + w_1 x_i) \right]^2 \quad \rightarrow \hat{y}$$

预测值 $\uparrow$

$$= \frac{1}{n} \sum \left[ (y_i - \bar{y}) - \underbrace{\frac{\sigma_{xy}}{\sigma_x^2}}_{w_1}(x_i - \bar{x}) \right]^2$$

$$= \frac{1}{n} \sum (y_i - \bar{y})^2 + \frac{1}{n} \sum \left( \frac{\sigma_{xy}}{\sigma_x^2} \right)^2 (x_i - \bar{x})^2 - 2 \cdot \frac{1}{n} \sum \frac{\sigma_{xy}(y_i - \bar{y})(x_i - \bar{x})}{\sigma_x^2}$$

$$= \sigma_y^2 + \frac{\sigma_{xy}^2}{\sigma_x^4} \cdot \sigma_x^2 - 2 \cdot \frac{\sigma_{xy}^2}{\sigma_x^2}$$

**FVU.**

$$= \sigma_y^2 + \frac{\sigma_{xy}^2}{\sigma_x^2} - 2 \frac{\sigma_{xy}^2}{\sigma_x^2}$$

$\dfrac{E(\hat{y} - y)^2}{Var(y)}$

Fraction of Variance Unexplained

(MSE=0)

$$MSE = \sigma_y^2 - \frac{\sigma_{xy}^2}{\sigma_x^2}$$

ideal $0 = \boxed{\dfrac{MSE}{\sigma_y^2}} = 1 - \boxed{\dfrac{\sigma_{xy}^2}{\sigma_x^2 \sigma_y^2}}$

ideal: $R^2 = 1$

$\boxed{R^2 \text{ value}}$  measure how good model is

$\dfrac{MSE}{\sigma_y^2} = 0 \quad (MSE = 0)$