$y_t$ \$

Scatter plot

4 years ago

today

rate $x_t$

$(r_t, p_t)$ $t=1,\ldots, T$

$y_t = \boxed{a_0} + b_1 x_t^{(1)} + \text{white noise } \boxed{\sigma} \varepsilon_t$

1

Goal: To find $a_0$ & $b_1$

$y_t = \boxed{a_0} + \boxed{b_1} x_t^{(1)} + \boxed{b_2} x_t^{(2)} + \boxed{b_3} x_t^{(3)} + \cdots$

$$y_t = \boxed{a_0} + \boxed{b_1} x_t^{(1)} + \varepsilon_t$$

1. linear regression (regress)

$a_0$ & $b_1$ & statistics

on $\boxed{R^2}$

$R^2 \leq 1$

2. Simplex

loss function

3. Gradient Descent (vanilla)

- $\boxed{\text{SGD \& Batch}}$

Y                    $\otimes$

↑                    ↑

observation          dependence

$[n \times 1]$

$[n \times d]$

$[n \times 1]$

$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad n \times 1$

$\otimes = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \begin{bmatrix} x_1^{(1)} \\ x_2^{(1)} \\ \vdots \\ x_n^{(1)} \end{bmatrix}$

$a_0 \qquad b_1$

$a_0, b_1, \ldots$

$X = -2 : 0.01 : 2 ; \quad \checkmark$

$X = -2 + 4 \times \text{rand} ; \quad \checkmark$

$\rightarrow \quad y_t = a_0 + b_1 x_t + \varepsilon_t$

$$\widehat{y_t} = a_0 + b_1 x_t \quad \leftarrow \text{model data}$$

$$y_t \qquad\qquad \leftarrow \text{true observation}$$

RMS

$$\sum_{t=1}^{T} \left( y_t - \widehat{y_t} \right)^2$$

$$= \frac{1}{T} \boxed{\sum_{t=1}^{T} \left( y_t - a_0 - b_1 x_t \right)^2} \quad \leftarrow \text{loss function}$$

Goal: To miminize the loss function distance for all t's

$$\ell(a_0, b_1) = \frac{1}{n} \sum_{j=1}^{n} \left( y_j - a_0 - b_1 x_j \right)^2$$

$$\nabla \ell = \begin{pmatrix} \frac{\partial \ell}{\partial a_0} \\ \frac{\partial \ell}{\partial b_1} \end{pmatrix} = \begin{pmatrix} -\frac{1}{n} \times 2 \sum_{j=1}^{n} \left( y_j - a_0 - b_1 x_j \right) \\ \frac{-2}{n} \sum_{j=1}^{n} \left( y_j - a_0 - b_1 x_j \right) x_j \end{pmatrix}$$

$$\left( \sum_{j=1}^{n} \cdots \right)$$

$\boxed{\text{Contour}}$ work up to 2 dimensional

$$+ \ 1$$

$\boxed{d > 2}$ $\qquad d = 5$

$\boxed{\text{loss function}}$ $\qquad\qquad \boxed{3\ 8}$

$$X_0 = [a_0, \ b_1, \ b_2, \ b_3, \ b_4]$$

$$X^* = [a_0^*, \ b_1^*, \ b_2^*, \ b_3^*, \ b_4^*]$$

$$X_{temp} = \alpha \, X_0 + (1-\alpha) X^*$$

$\alpha = 1$ $\qquad\qquad \boxed{X_{temp} = X_0}$

$\alpha = 0$ $\qquad\qquad \boxed{X_{temp} = X^*}$

loss function

$\alpha$

$-1$    $2$

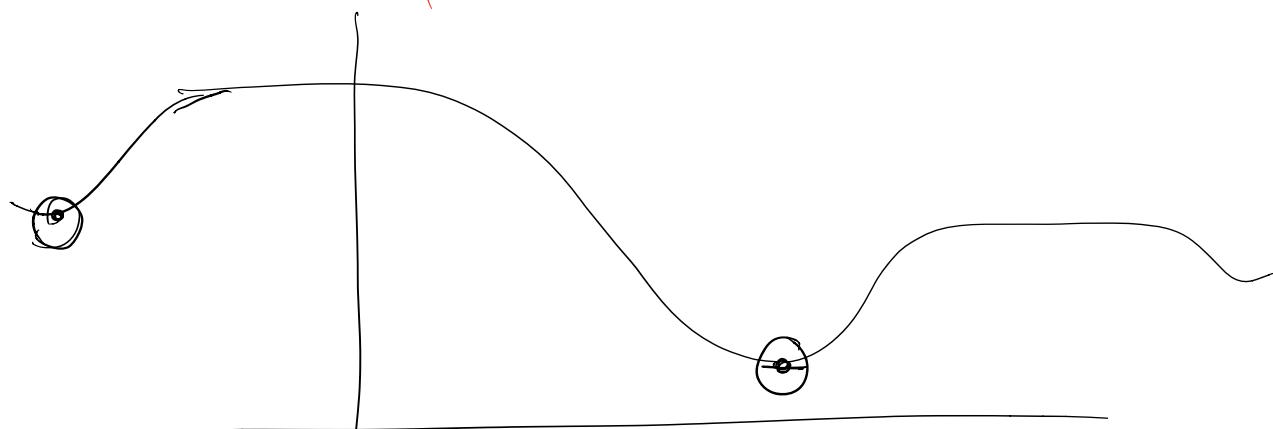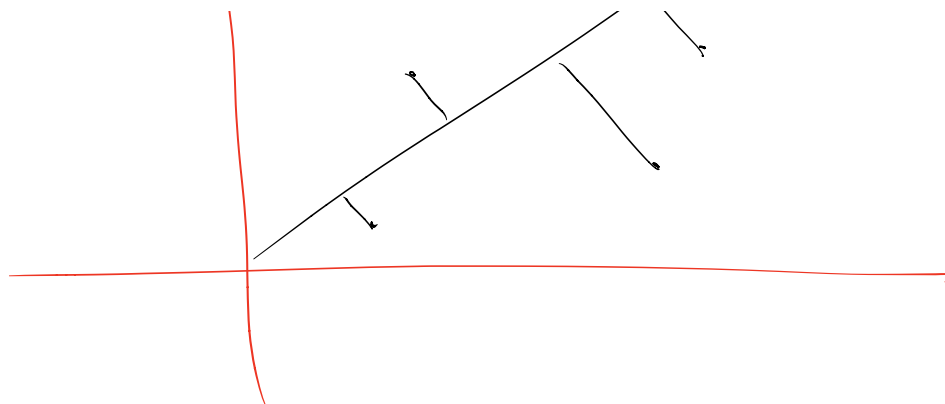$X_0^{(1)}$    $X^*$

$X_0^{(2)}$

Vanilla GD
SGD
BGD

- Learning Algorithm

      ✓. Statistical ML / Bayesian

                              updates

      ✓. ML

      ✓. DL

Loss functions & Optimization
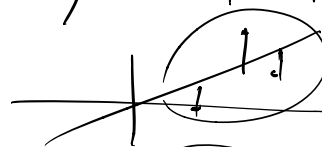
Loss functions
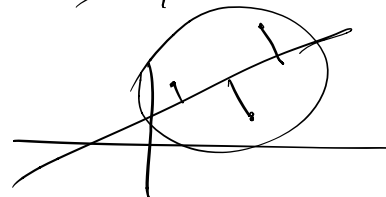
- Choice of starting point

- Choice of optimization

      Simplex vs. Gradient Descent

- Choice of objective function / loss function

      - distance

      - y difference

- Stress-testing ✓