# Lecture November 22

$$\text{Accuracy score} = \frac{\sum_{i=0}^{m-1} I(y_i' = \tilde{y}_i)}{m}$$

confusion matrix

TP = True Positive
equ with a correct
classification

TN = True negative, equ
with correct rejection

FP = False positive
false alarm

FN = False negative, equ
with miss

True positive rate

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

False positive rate

$$= \frac{FP}{FP + TN} \neq FPR$$

True Negative rate

$$= TNR = \frac{TN}{TN + FP}$$

Gains curve

$$\frac{count\ TP + count\ FP}{all\ observations}$$

$$FPR = 1 - TNR$$

# Gradient descent

$n = $ # data points

$P = $ # features

$n \gg p \quad \sim \quad n \gtrsim p \quad n \leq 10^5$

suppose we have the optimal
$\hat{\beta}$, in principle this is an
iterative process

$$\hat{\beta} \cong \beta_{k+1}$$

$$|\beta_{k+1} - \beta_k| \leq \varepsilon \approx 10^{-10}$$

$$\beta = \beta_k, \quad \text{Taylor expand}$$

$$C(\hat{\beta}) = C(\beta) + g^T(\hat{\beta} - \beta)$$

$$\nabla_\beta C(\beta)$$
$$\nearrow$$
$$\text{evaluated at}$$
$$\beta$$

$$+ \frac{1}{2}(\hat{\beta} - \beta)^T H (\hat{\beta} - \beta)$$

$$\nearrow$$
$$\frac{\partial^2 C(\beta)}{\partial \beta \, \partial \beta^T}$$

$$\nearrow$$
$$\text{Logistic reg}$$
$$x^T W x$$

Define $\quad b = \hat{\beta} - \beta$

$$C(\hat{\beta}) = C(\beta) + b^T g +$$
$$\frac{1}{2} b^T H b$$

$$\frac{\partial C}{\partial b^T} = 0 = Hb + g \quad \Rightarrow$$

$$b = \hat{\beta} - \beta = -H^{-1} g \quad \Rightarrow$$

$$\hat{\beta} = \beta - H^{-1}g$$

$$\beta_{k+1} = \beta_k - H^{-1}(\beta_k) \, \nabla_\beta C(\beta_k)$$

Newton- Raphsen's method.

Newton's method is derived from a general function

$$f(x) = \frac{1}{2} x^T A x + x^T b + c$$

$$\frac{\partial f(x)}{\partial x^T} = 0 = Ax + b = 0 \Rightarrow$$

$$Ax = -b = 0$$

Known      Known

algorithm:
- start with guess $\beta_0$
- iterate till

$$|\beta_{k+1} - \beta_k| \leq \varepsilon$$

$$\beta_{k+1} = \beta_k - H^{-1}(\beta_k) \, \nabla_\beta C(\beta_k)$$

since we have to compute $H^{-1}$ repeatedly, replace $H^{-1}(\beta_k)$ with a constant $\gamma_k$ $\Rightarrow$

$$\beta_{k+1} = \beta_k - \gamma_k D_\beta C(\beta_k)$$

$$\equiv \text{GRADIENT DESCENT}$$

$$b = \hat{\beta} - \beta = -H^{-1} \cdot g =$$
$$-H^{-1} D_\beta C(\beta)$$

$H^{-1} \Rightarrow$ learning rate $\gamma$

$$\hat{\beta} - \beta = -\gamma D_\beta C(\beta) = -\gamma g(\beta)$$

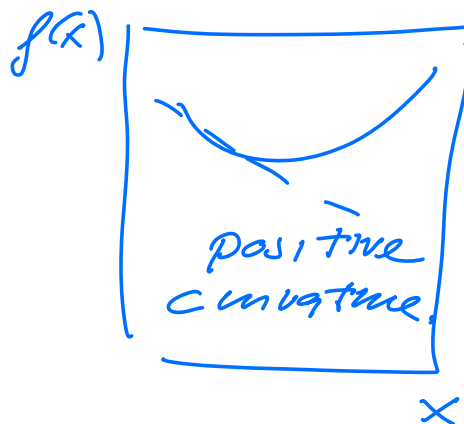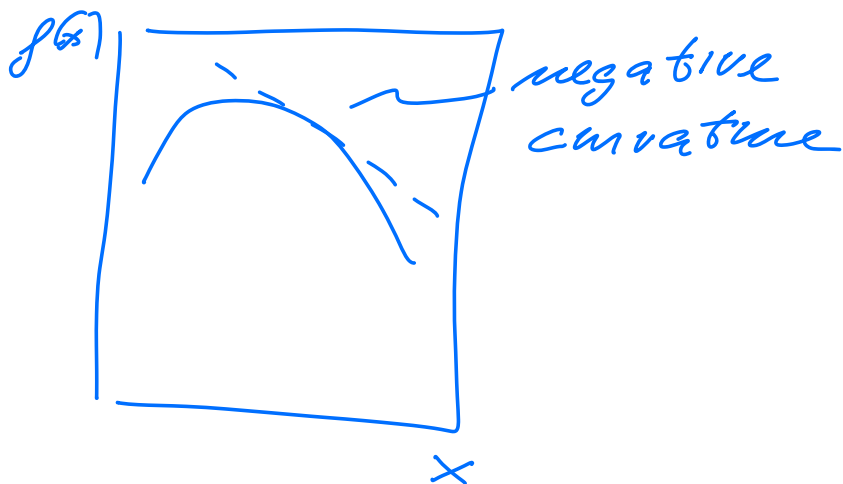$$C(\hat{\beta}) = C(\beta) - \gamma g^T g$$
$$+ \frac{1}{2} \gamma^2 g^T H g$$

in one-dimension

$$C(\vec{\beta}) = C(\beta) - \gamma g^2 + \frac{1}{2}\gamma g^2 H$$

$C(\beta) = $ original $C$ / start value

$\gamma g^2 = $ improvement to the slope of $C(\vec{\beta})$

$\frac{1}{2}\gamma^2 g^2 H = $ correction due to curvature



$f(x)$    negative curvature

   no curvature

$f(x)$    positive curvature

$x$

$$\partial C \qquad \qquad \qquad ^2 \qquad g^2 H$$

$$\frac{\partial}{\partial \gamma} = 0 = -g + \gamma g''$$

or $\qquad -g^T g + \gamma g^T H g =>$

$$\gamma = \frac{g^T g}{g^T H g} \left( = \frac{g^T g}{g^T g \lambda} = \frac{1}{\lambda} \right)$$

if $\quad H g = \lambda H$

$$\gamma = \frac{1}{\lambda}$$

smallest $\gamma = 1/\lambda_{max}$

largest $\quad \gamma = 1/\lambda_{min}$.

For convergence of Newton-Raphsen we must have

$$\gamma < \frac{2}{\lambda_{max}}$$

where $\lambda_{max}$ is the largest eigenvalue of $H$.

$f(x)$    $x_0$    **Too small**    $\gamma \geq \frac{2}{\lambda_{max}}$

$x_2$

$x_1$

$\gamma < \frac{2}{\lambda_{max}}$

(if not too small)

$x_0$

$x$

$x^{opt} = \hat{x}$

$f(x)$

saddle point

$x_0$

ideal global min

local min

$x$