$$\hat{G}_{m-1}(x) = \sum_{i=1}^{m-1} \beta_i \hat{f}_i(x)$$

Solving $\beta_m, \hat{f}_m = \underset{\beta, f_m}{\arg\min} \sum_n L\left(y_n, \hat{G}_{m-1}(x_n) + \beta f_m(x_n)\right)$

$\hookrightarrow L(y, h) = e^{-yh}$

$$\sum_n \exp\left(-y_n\left(\hat{G}_{m-1}(x) + \beta f_m(x_n)\right)\right)$$

$$= \sum_n \exp\left(-y_n \hat{G}_{m-1}(x)\right) \exp\left(-y_n \beta f_m(x_n)\right)$$

$\underbrace{\qquad\qquad\qquad\qquad}$
no $\beta, f_m$

So I can regard as a weight $\omega_{nm}$

$$= \sum_n \omega_{nm} \exp\left(-\beta \underbrace{y_n f_m(x_n)}_{\pm 1}\right)$$

$y_n = \pm 1$

$f_m(x_n) = \pm 1$

$\underbrace{\qquad\qquad\qquad}_{e^{\pm\beta}}$

$$= \sum_{y_n = f_m(x_n)} \omega_{nm} e^{-\beta} + \sum_{y_n \neq f_m(x_n)} \omega_{nm} e^{\beta}$$

$$= \underbrace{e^{-\beta} \sum_{y_n = f_m(x_n)} \omega_{nm}}_{①} - \underbrace{e^{-\beta} \sum_{y_n \neq f_m(x_n)} \omega_{nm}}_{③} + \underbrace{e^{-\beta} \sum_{y_n \neq f_m(x_n)} \omega_{nm}}_{②} + \underbrace{e^{\beta} \sum_{y_n \neq f_m(x_n)} \omega_{nm}}_{④}$$

$$= \boxed{① + ②} \quad + \boxed{③ + ④}$$

$$= e^{-\beta} \sum_n \omega_{nm} + \left(e^{\beta} - e^{-\beta}\right) \sum_{y_n \neq f_m(x_n)} \omega_{nm}$$

$$⊛ = e^{-\beta} \sum_n \omega_{nm} + \left(e^{\beta} - e^{-\beta}\right) \sum_n \omega_{nm} \mathbb{1}(y_n \neq f_m(x_n))$$

$\underbrace{\qquad\qquad\qquad\qquad}$
$\propto$ weighted
mis-class err of $f_m$

Want to find opt. $\beta, f_m$ to minimize

- Fix $\beta$, find $f_m$ to minimize:

  $\hat{f}_m$ should minimize weighted mis-class err rate.

Given $\hat{f}_m$, find best $\beta$:

$$\frac{\partial}{\partial \beta}(⊛) = -e^{-\beta} \sum_n \omega_n + \left(e^{\beta} + e^{-\beta}\right) \sum_n \omega_{nm} \mathbb{1}(y_n \neq \hat{f}_m(x_n))$$

$$= 0 \qquad \text{and solve for } \beta$$

$$\Leftrightarrow \sum_n \omega_n = \left(e^{2\beta} + 1\right) \sum_n \omega_{nm} \mathbb{1}(y_n \neq \hat{f}_m(x_n))$$

$$\vdots$$

$$\beta_m = \frac{1}{2} \log\left(\frac{1 - err_m}{err_m}\right)$$

$err_m =$ weighted mis-class

$$\alpha_m \triangleq 2\beta_m$$

Just need to show that we update weights as claimed.

$$\boxed{\text{Ada Boost}: \quad w_{n\,m+1} \leftarrow w_{nm}\,\exp\!\left(\alpha_m\,\mathbb{1}(y_n \neq \hat{f}_m(x_n))\right)}$$

Notice:

$$
\begin{aligned}
w_{n,m+1} &= \exp\!\left(-y_n\,\hat{G}_m(x_n)\right)\\[4pt]
&= \exp\!\left(-y_n\,\hat{G}_{m-1}(x_n)\right)\exp\!\left(-\beta_m\,\hat{f}_m(x_n)\,y_n\right)\\[4pt]
&= w_{nm}\,\exp\!\left(-\beta_m\,\hat{f}_m(x_n)\,y_n\right)
\end{aligned}
$$

notice that $\;\underbrace{-y_n\,\hat{f}_m(x_n)}_{\pm 1} = 2\,\mathbb{1}(y_n \neq \hat{f}_m(x_n)) - 1$

So

$$
\begin{aligned}
&= w_{nm}\,\exp\!\left(2\beta_m\,\mathbb{1}(y_n \neq \hat{f}_m(x_n)) - \beta_m\right)\\[4pt]
&= w_{nm}\,\exp\!\left(\alpha_m\,\mathbb{1}(y_n \neq \hat{f}_m(x_n))\right)\exp\!\left(-\beta_m\right)
\end{aligned}
$$

doesn't matter

This is exactly Ada Boost.

For continuos $y_n$, can do a similar thing, by fitting each model to the residual of the previous.

This is equivalent to SAM using a Squared err loss.

---

A couple tuning choices:

- ~ loss *
- — number of trees $M$
- — Shrinkage factor: $v \in [0,1]$

$$\hat{G}_m(x) = \hat{G}_{m-1}(x) + v \, \alpha_m \hat{f}_m(x)$$

<span style="color:red">↑ learning rate</span>

- — number of splits in each of my trees

---

Gradient Boosting

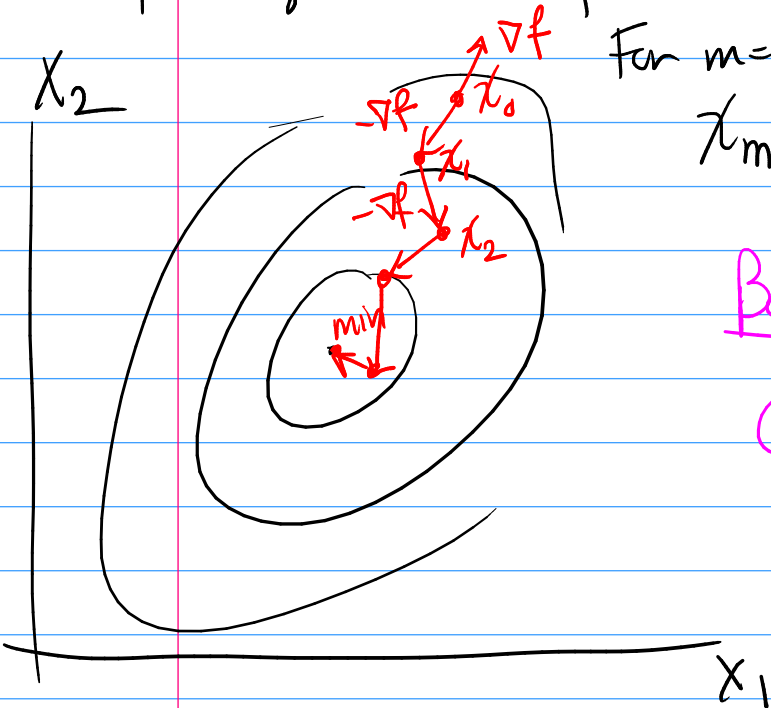Boosting using something other than SE loss (regr.) or exp. loss (class.)

# Gradient Descent:

(minimize)

Optimize some function $f$

$X_2$



For $m = 1, \ldots, M$
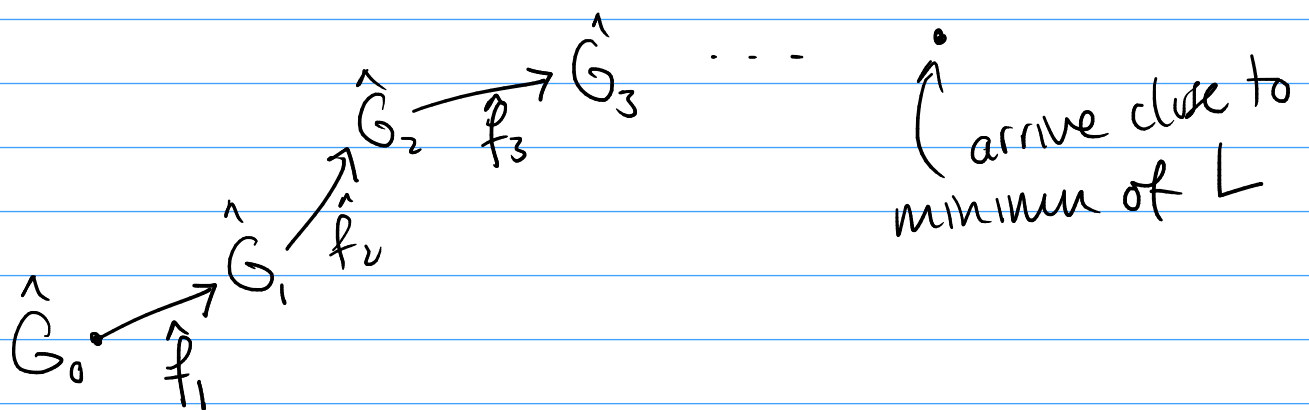
$$X_m = X_{m-1} - \alpha_m \nabla f \big|_{X_{m-1}}$$

## Boosting:

$$\hat{G}_m(x) = \hat{G}_{m-1}(x) + \alpha_m \hat{f}_m(x)$$

Looks like gradient descent

if $\hat{f}_m \approx -\nabla f = -$ gradient

$$\hat{G}_0 \bullet \xrightarrow{\hat{f}_1} \hat{G}_1 \xrightarrow{\hat{f}_2} \hat{G}_2 \xrightarrow{\hat{f}_3} \hat{G}_3 \cdots \uparrow \bullet$$

arrive close to minimum of $L$

## Gradient Boosting : <span style="color:red">(regression)</span>

(0) $\hat{G}_0(x) = 0$

(2) For $m = 1, \ldots, M$

(a) $r_{nm} = - \left. \frac{\partial L}{\partial f(x_n)} \right|_{\hat{G}_{m-1}(x_n)}$     pseudo residuals

(b) fit $\hat{f}_m$ to predict $r_{nm}$ from $x_n$

(c) $\hat{G}_m(x) = \hat{G}_{m-1}(x) + \alpha_m \hat{f}_m(x)$

<span style="color:magenta">learning rate</span>