# MATH 565 : Lecture 8 (02/05/2026)

Today : 
* SGD for sum
* Logistic regression loss
* coordinate descent (CD)

---

**Recall**

$$\nabla J_{H\text{-}SVM} = -y_i \bar{x}_i \, \delta\left(1 - y_i(\bar{w}^T \bar{x}_i) > 0\right) + \lambda \bar{w} \qquad \text{(Hinge-SVM)}$$

$\hookrightarrow$ indicator function

$$\nabla J_{L_2\text{-}sum} = -y_i \bar{x}_i \max\left\{0, \, 1 - y_i(\bar{w}^T \bar{x}_i)\right\} + \lambda \bar{w} \qquad (L_2\text{-SVM})$$

## SGD for H-SVM

$$\bar{w} \leftarrow \bar{w}(1 - \alpha\lambda) + \alpha \sum_{\substack{i=1 \\ i \in S}}^{n} y_i \bar{x}_i \, \delta\left(1 - y_i(\bar{w}^T \bar{x}_i) > 0\right)$$

Let $S^+ = \left\{ i \in S \,\middle|\, y_i(\bar{w}^T \bar{x}_i) < 1 \right\} \longrightarrow$ subset of indices in $S$ for which $\delta(\cdot) = 1$ above

$y_i(\bar{w}^T \bar{x}_i) < 0 \implies i$: misclassified point/instance

$y_i(\bar{w}^T \bar{x}_i) \in (0,1) \implies i$: correctly classified instance, but lies close to decision boundary.

When $y_i(\bar{w}^T \bar{x}_i) \geqslant 1$, $i$ is correctly classified and well-separated $\implies$ does not contribute to $J$.

## SGD for Hinge-SVM

$$\bar{w} \leftarrow \bar{w}(1 - \alpha\lambda) + \sum_{i \in S^+} \alpha y_i \bar{x}_i \longrightarrow \text{primal SVM algorithm}$$

proposed by Hinton in 1989 !

$\hookrightarrow$ before VC dimension and other details were proposed by Vapnik and coauthors later on...

# Logistic Regression Loss

Note that the hinge loss function, while convex, is not smooth — there is a sharp "hinge" at the value of 1 (hence the name). The logistic regression loss can be considered as a smooth version of the hinge loss.

$$J_{LR} = \sum_{i=1}^{n} \log\left(1 + e^{-y_i(\bar{w}^T\bar{x}_i)}\right) + \frac{1}{2}\|\bar{w}\|^2$$

Consider $L(z) = \log\left(1 + e^{-z}\right)$ with $z = yf(\bar{x})$ as the prediction

$$= \log\left(e^{-z}(1 + e^{z})\right)$$

$$= -z + \underbrace{\log\left(1 + e^{z}\right)}_{\longrightarrow 0 \text{ as } z \longrightarrow -\infty}$$

e.g., $z_i = y_i(\bar{w}^T\bar{x}_i)$
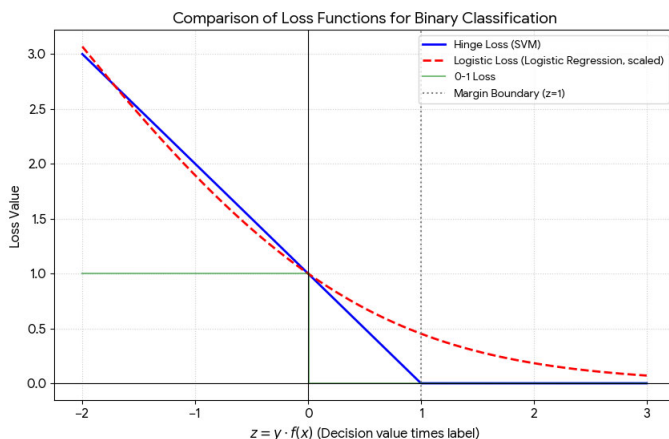
largely misclassified instances give huge negative values

$\Rightarrow$ For largely misclassified instances, $J_{LR}$ increases linearly as $|\bar{w}^T\bar{x}_i|$ increases. For such instances

<span style="color:red">recall that the hinge loss function is $1-z$ for $z < 1 \ldots$</span>

$$J_{H\text{-}SUM} - J_{LR} \approx 1$$

$\Rightarrow$ H-SUM and LR-SVM treat grossly misclassified instances similarly.

But $J_{LR} > 0$ (ignoring $\|\bar{w}\|^2$ term) for all instances.



Comparison of Loss Functions for Binary Classification

- Hinge Loss (SVM)
- Logistic Loss (Logistic Regression, scaled)
- 0-1 Loss
- Margin Boundary (z=1)

$z = y \cdot f(x)$ (Decision value times label)

Loss Value

$J_{LR}$ is differentiable while $J_{H\text{-}SUM}$ is not.

In fact, the fit term in the $J_{LR}$ loss function turns out to be strictly convex (on its own)!

**Lemma 8** $J_{LR}$ without the $\frac{1}{2}\|\bar{w}\|^2$ regularizer term is strictly convex.

But we usually add $\frac{1}{2}\|\bar{w}\|^2$ still, as this extra term encourages sparsity. We could instead use $\sum_{i=1}^{n}|w_i|$ as an $L_1$-regularity term. But then again, $|w_i|$ is not smooth either...

We finish with the details of gradient descent using the logistic regression loss function for SVM.

$$\nabla J_{LR} = -\sum_{i=1}^{n}\frac{y_i \bar{x}_i}{\left[1+e^{-y_i(\bar{w})\bar{x}_i)}\right]} + \lambda\bar{w}$$

Hence, the SGD update is given as follows.

$$\bar{w} \leftarrow \bar{w}(1-\alpha\lambda) + \sum_{i\in S}\frac{\alpha \, y_i \bar{x}_i}{\left[1+e^{-y_i(\bar{w}^T\bar{x}_i)}\right]}$$

# Coordinate Descent (CD)

Recall the gradient descent update: $\bar{w} \leftarrow \bar{w} - \alpha \nabla J$.

In coordinate descent (CD), we optimize one coordinate at a time.

$$\bar{w} = \arg\min_{\bar{w}} \left\{ J(\bar{w}) \mid \text{only } w_i \text{ varies} \right\}$$

\* only one variable to handle — can be (much) easier.
\* can use line search if not able to solve exactly.

Cycle through all $i = 1, \ldots, d$. → in a full cycle
if no $w_i$ changes, <u>STOP</u>.

\* If J is convex and differentiable, then the converged solution is optimal.
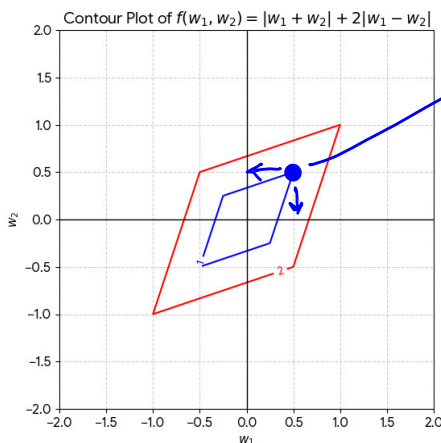\* But if J is <u>not</u> differentiable, this guarantee does not hold, even when it may be convex.

Consider $J(\bar{w}) = |w_1 + w_2) + \alpha |w_1 - w_2|$, $\alpha > 1$. → can show that J is convex.

J is minimal at $(0,0)$.

two contours of J=1, J=2.



Contour Plot of $f(w_1, w_2) = |w_1 + w_2| + 2|w_1 - w_2|$

But if we're at $(1,1)$, neither coordinate will decrease J.

Still, we can give a characterization of a fairly general class of loss functions on which CD is well-behaved.

> single variable functions!

**Lemma 9**   Let $J(\bar{w}) = G(\bar{w}) + \sum_{i=1}^{d} H_i(w_i)$   where $G(\bar{w})$ is convex and differentiable, while $H_i(w_i)$ are convex but may not be differentiable. Then, CD converges to the global minimum of $J$.

The $L_1$-regularizer has this structure: $H_i = |w_i|$, giving

$$\sum_{i=1}^{d} |w_i| \text{ as the regularizer term.}$$

But in some cases, variable transformations can help even if $J$ is not in this form.

e.g., $J(\bar{w}) = G(\bar{w}) + |w_1 + w_2| + \alpha|w_1 - w_2|$,   $\alpha > 1$. Here,

$\hookrightarrow$ convex, differentiable

we can use $u_1 = \dfrac{w_1 + w_2}{2}$ and $u_2 = \dfrac{w_1 - w_2}{\alpha}$ to get

$$w_1 = u_1 + u_2 \text{ and } w_2 = u_1 - u_2, \text{ giving}$$

$$J(\bar{u}) = G(u_1 + u_2, u_1 - u_2) + 2|u_1| + 2\alpha|u_2|,$$

which has structure specified in Lemma 9.

$\hookrightarrow$ CD will work well on this version of $J$.

# Linear Regression with Coordinate Descent

To understand CD better, we apply it to linear regression (without regularization)

$$J(\bar{w}) = \frac{1}{2}\|D\bar{w} - \bar{y}\|^2 = \frac{1}{2}\sum_{i=1}^{n}(\bar{w}^T \bar{x}_i - y_i)^2$$

$$= \frac{1}{2}\sum_{i=1}^{n}\left(\underbrace{\sum_{j=1}^{d}w_j x_{ij} - y_i}\right)^2 \qquad \text{We consider CD for } w_k$$

$$\longrightarrow w_k x_{ik} + \sum_{j \neq k}w_j x_{ij} - y_i$$

$$\frac{\partial J}{\partial w_k} = \sum_{i=1}^{n}\left(w_k x_{ik} + \sum_{j \neq k}w_j x_{ij} - y_i\right)x_{ik} = 0 \quad \longrightarrow \text{first order optimality}$$

$$\longrightarrow \left(y_i - \sum_{j \neq k}w_j x_{ij}\right)$$

$$\Rightarrow \quad w_k = \frac{-\sum_{i=1}^{n}\left(\sum_{j \neq k}w_j x_{ij} - y_i\right)x_{ik}}{\sum_{i=1}^{n}x_{ik}^2} \qquad \begin{array}{l}\text{term}\\\text{not included}\end{array}$$

With $\bar{r} = \bar{y} - D\bar{w} = \bar{y} - \sum_{j \neq k}\bar{d}_j w_j + \boxed{w_k \bar{d}_k}$ $\qquad \bar{d}_j : j^{th}$ column of $D$.

$\hookrightarrow$ vector of residuals

The update step is gives as follows.

$$w_k^{new} = \frac{\bar{d}_k^T\left(\bar{r} + w_k^{old}\bar{d}_k\right)}{\|\bar{d}_k\|^2}$$

$$D = \begin{array}{c}\begin{array}{cccccc}1 & 2 & \cdots & k & & d\end{array}\\ \begin{array}{c}1\\2\\\vdots\\\\n\end{array}\left[\begin{array}{c}\boxed{\begin{array}{c}x_{1k}\\x_{2k}\\\vdots\\x_{ik}\\\vdots\\x_{nk}\end{array}}\end{array}\right]\end{array}$$

$$= w_k^{old} + \frac{\bar{d}_k^T\bar{r}}{\|\bar{d}_k\|^2} \qquad \longrightarrow \begin{array}{l}\text{we assume the trivial case}\\\text{of } \bar{d}_k = \bar{0} \text{ does not occur.}\end{array}$$

If data columns are normalized, $\|\bar{d}_k\|^2 = 1$, and we get $w_k^{new} = w_k^{old} + \bar{d}_k^T\bar{r}$.

More generally, we update

$$w_k^{new} \leftarrow w_k^{old} + \bar{d}_k^T \bar{r} \quad \rightarrow \quad \text{extremely efficient!}$$

We also update the residuals $\bar{r}$ as follows:

$$\bar{r} \leftarrow \bar{r} - \bar{d}_k \overbrace{(\Delta w_k)}^{w_k^{new} - w_k^{old}}$$