

Last time

- Abridged history of deep learning
- Charted the convergence of the key innovations that led to AlexNet in 2012

△ only 13 years ago

- model development (F)

- backpropagation (opt)

- GPUs (opt)

- ImageNet (D)

Today:

- statistical learning framework

- F (search space)

- L (search criterion)

- opt (search algorithm)

- gradient descent

Def Statistical Learning

→ a principled process, informed by data, for the optimal selection of a function (\hat{f}) to accomplish a task (T).

Def Deep Learning

→ statistical learning with the set of candidate functions (F) being neural networks

Statistical Learning Process

Step 1: Selection of Search Space (F)

→ First, we need to define which functions (f) we are choosing between.

→ Restrict the search to a certain set of functions $f \in F$

e.g. $T = \text{regression} : F = \{ \beta_0 + \beta_1 X : \beta_0, \beta_1 \in \mathbb{R} \}$

e.g. $T = \text{inference} : F = \left\{ \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2}\left(\frac{x_i - \mu}{\sigma}\right)^2\right\} : \begin{matrix} \mu \in \mathbb{R}, \sigma^2 \in \mathbb{R}^+ \end{matrix} \right\}$

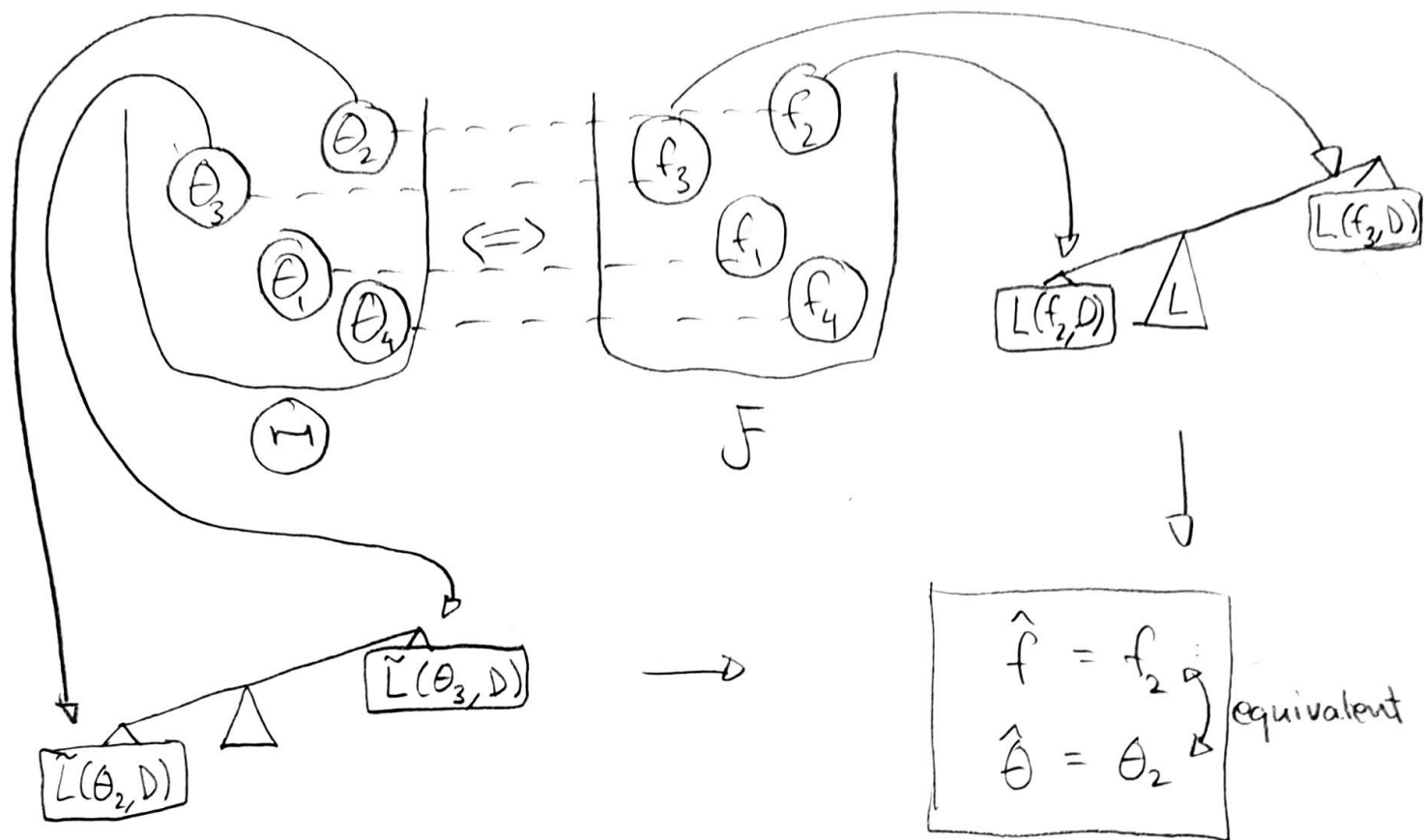
Step 2: Selection of Search Criterion (L)

- want to choose the best f among $f \in \mathcal{F}$
- need a way to evaluate task performance of each $f \in \mathcal{F}$
- Define a loss function $L(f, D)$ that scores a given $f \in \mathcal{F}$ on how (badly) it accomplishes the task on the observed data (D).
- Then search for the best (as measured by L) $f \in \mathcal{F}$. Ideally obtaining:

$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{F}} L(f, D)$$

$$\Leftrightarrow f_{\hat{\theta}} = \operatorname{argmin}_{\theta \in \Theta} L(f_{\theta}, D)$$

$$\Leftrightarrow \hat{\theta} = \operatorname{argmin}_{\theta \in \Theta} \tilde{L}(\theta, D)$$



Challenges

- ① How to encode performance on task T into a loss function?
- ② How to conduct the search of F ? (i.e. optimization)

① Loss functions (the creative bit)

→ Task: regression

→ mean squared error : $L(f, D) = \frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2$

→ mean zero-one loss : $L(f, D) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{f(x_i) \neq y_i\}$

→ Task: classification

→ cross entropy loss : $L(f, D) = \frac{1}{n} \sum_{i=1}^n - \sum_{c=1}^C \log(f_c(x_i)) \mathbb{I}\{y_i = c\}$

→ misclassification rate : $L(f, D) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{\arg\max_{c \in \{1, \dots, C\}} f_c(x_i) \neq y_i\}$

Different choices for L will affect how each candidate is scored and may thus determine different models to be optimal.

② Optimization

Our task now becomes solving:

$$\hat{\theta} = \underset{\theta \in \mathcal{H}}{\operatorname{argmin}} L(\theta, D)$$

using some algorithm, call it "opt", so

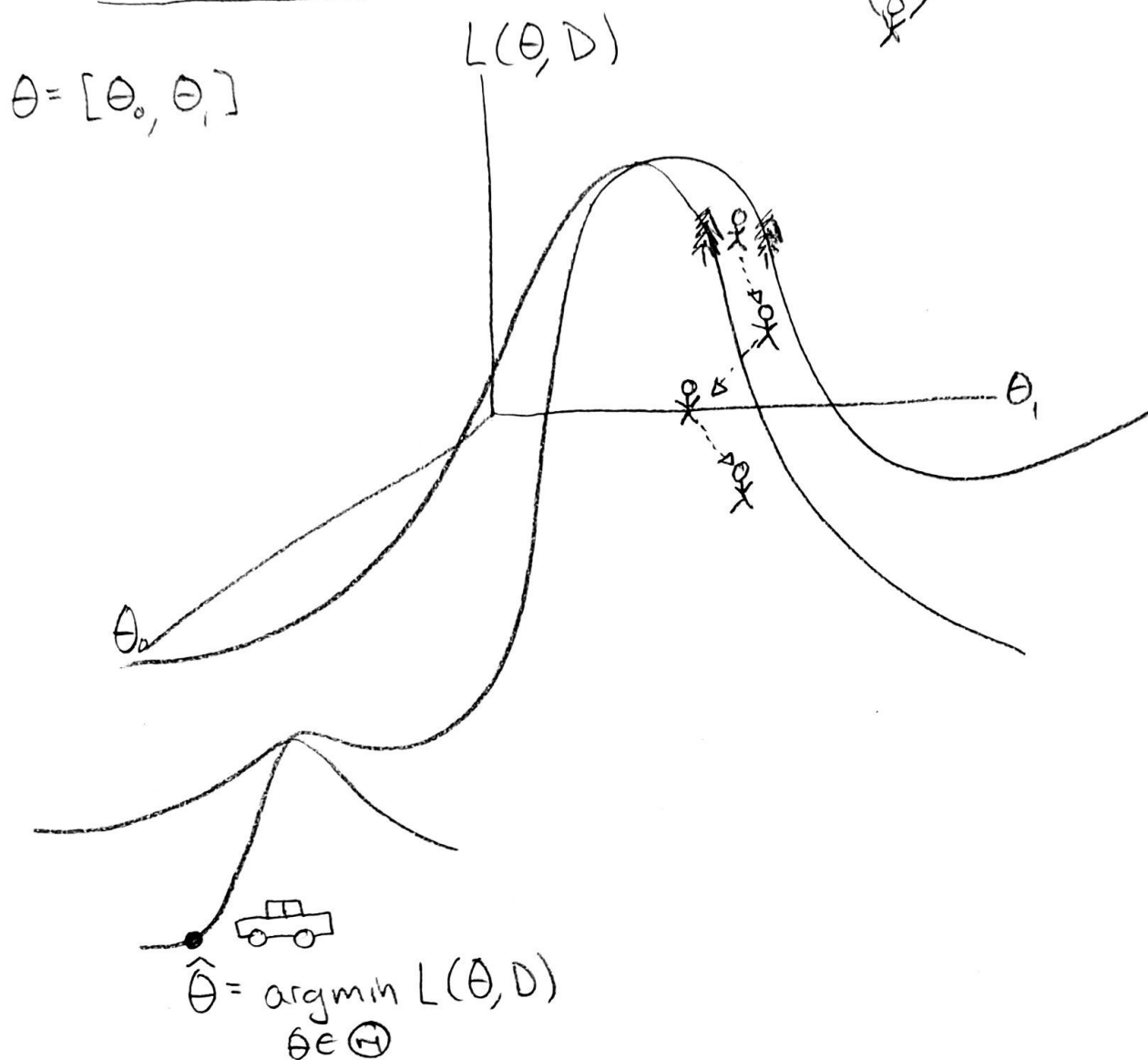
$$\hat{\hat{\theta}} = \underset{\theta \in \mathcal{H}}{\operatorname{opt}} L(\theta, D)$$

where we would like $|L(\hat{\theta}, D) - L(\hat{\hat{\theta}}, D)|$ to be small.

Some approaches:

- Don't search Θ , just pick a $\theta \in \Theta$ at random
- Exhaustive search: evaluate $L(\theta, D) \forall \theta \in \Theta$
- Analytical solution (calculus)
- Random Search: evaluate $L(\theta, D) \forall \theta \in \Theta, \subset \Theta$
- Gradient Descent & Variants
 - Stochastic Gradient Descent (SGD)
 - SGD + Momentum
 - RMS prop
 - Adam (SGD + momentum + RMS prop)

Gradient Descent



Idea: Imagine the loss function as a surface. You are standing on the surface and want to walk to the lowest point, but you can't see very far. Look down, determine in which direction the ground is sloped most steeply (the gradient), walk a step in this direction, then repeat.

Algorithm: given the function $L(\theta)$, with Jacobian (ie multivariate derivative) $\nabla_{\theta} L(\theta) = \frac{\partial L(\theta)}{\partial \theta}$, then $\theta^{(t+1)}$, the $(t+1)^{\text{th}}$ estimate of $\underset{\theta \in \Theta}{\operatorname{argmin}} L(\theta)$,

is given as follows:

$$\theta^{(t+1)} = \theta^{(t)} - \alpha_t [\nabla_{\theta} L(\theta)] \Big|_{\theta = \theta^{(t)}}, \quad \forall t = 0, \dots, T-1$$

where $\alpha_t > 0$ is the "step size" or "learning rate"

→ $\theta^{(0)}$ can be initialized randomly or using prior knowledge

→ T , the number of steps, can be chosen beforehand or can be a function of the algorithm's progress (ie a stopping rule)