

Deep learning.

representation learning; end-to-end learning → params in each layer are trained w.r.t. overall loss.

input → learned feature extractor → learned features → learned classifier → label.

Big Model; Large Datasets; Enough Compute.

good at scaling → the ability to work better w/ more data

inductive bias VS learning.

Machine Learning Review:

① Supervised Learning. Model → Loss fn → optimizer.

- output probability instead of output themselves. / probabilistic models

$$\text{Softmax}(f_{\theta}(x))_c = \frac{\exp(f_{\theta}(x)_c)}{\sum_i \exp(f_{\theta}(x)_i)} = \text{Poly}(c|x) \quad k \text{ possible labels.}$$

↑ entire distribution over all logits
make output positive

- Maximum Likelihood Estimation. MLE * consistency.

$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}. \quad \Omega_{\text{MLE}} = \underset{\theta \in \Theta}{\operatorname{argmax}} P(D|\theta) = \underset{\theta \in \Theta}{\operatorname{argmax}} \prod_{i=1}^N P(x_i) P(y_i|x_i)$$

Negative Log Likelihood:

$$\underset{\theta \in \Theta}{\operatorname{argmin}} \sum_{i=1}^N -\log P(x_i) + \log P(y_i|x_i) = \underset{\theta \in \Theta}{\operatorname{argmax}} \sum_{i=1}^N \log P(y_i|x_i) \Rightarrow \ell(\theta; x_i, y_i) = \underset{\theta \in \Theta}{\operatorname{argmin}} \sum_{i=1}^N -\log \text{Poly}(y_i|x_i)$$

Cross Entropy Loss: (Maximizing log likelihood ≈ minimizing cross-entropy).

$$H(p, q) = -\sum_x p(x) \log q(x) = E_p[-\log q(x)].$$

$$H(p_{\text{data}}, p_{\theta}) = E_{\text{data}}[-\log p_{\theta}(x|y)] = E_{\text{data}}[-\log p(x) - \log p_{\theta}(y|x)] \approx \sum_{i=1}^N -\log(p_{\theta}(x_i))$$

↑ const. w.r.t. θ.

- Iterative optimization. / gradient based optimization.

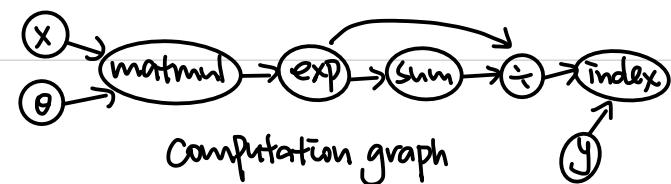
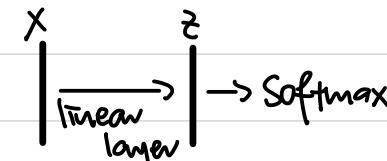
$$\theta \leftarrow \theta - 2 \Delta \theta \frac{1}{N} \sum_{i=1}^N \ell(\theta; x_i, y_i)$$

0-1 loss fn is not differentiable.

- Logistic regression

$$\ell(\theta; x, y) = -\log P_{\theta}(y|x) = -\log \text{Softmax}(f_{\theta}(x))_c.$$

$$\theta \leftarrow \theta - 2 \Delta \theta \frac{1}{N} \sum_{i=1}^N \ell(\theta; x_i, y_i)$$



- θ COULD be a $d \times (k-1)$ matrix.

Example: $k=2$ $P(y=1|x) = \text{Softmax}([\theta_0, \theta_1]^\top [x])_1 = \frac{\exp^T x}{\exp^T x + 1}$

$$\ell(\theta; x_i, y_i) = (1-y_i)(\log(\exp^T x + 1)) - y_i(\theta^T x - \log(\exp^T x + 1)).$$

$$\nabla \ell(\theta; x_i, y_i) = \left(\frac{\exp^T x}{\exp^T x + 1} - y_i \right) x.$$

- true risk $R(\theta) = \mathbb{E}[\ell(\theta; x, y)]$ whole dataset.

empirical risk $\hat{R}(\theta) = \frac{1}{N} \sum_{i=1}^N \ell(\theta; x_i, y_i)$ ERM. training set.

$R(\theta) \neq \hat{R}(\theta)$ can't reuse the training data.

high true risk, low empirical risk. \rightarrow overfitting. too small dataset / too powerful model

high true risk, high empirical risk \rightarrow underfitting. weak model / bad optimization.

In general, true risk \geq empirical risk.

- model class: the set of all possible fns that the model can represent via different params.

capacity = |model class|.

training set - learning θ . test underfitting. optimization.

Validation set - estimating $R(\theta)$. test overfitting.

test set - only use once.

combat overfitting, even underfitting $f_2 \text{ norm} = \text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$ $f_1 \text{ norm} = \text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$

- regularizer: anything we add to the loss fn optimization that not depend on data.

Maximum a posteriori (MAP.) $\xrightarrow{\text{inductive bias}}$

$$\sum_{i=1}^n -\log p(y_i|x_i, \theta) - \log p(\theta). \quad p(\theta) = N(\theta; 0, \sigma^2 I), \text{ so } -\log p(\theta) = \sum_{i=1}^n \frac{1}{2} \frac{\theta_i^2}{\sigma^2} + \text{Const.} = \lambda \|\theta\|_2^2 \text{ where } \lambda = \frac{1}{2\sigma^2}$$

$$\sum_{i=1}^n -\log p(y_i|x_i, \theta) + \lambda \|\theta\|_2^2 \quad | \quad f_2 \text{ regularization.} = \text{Ridge regression.}$$

hyperparameter, "weight decay". Smaller $\lambda \rightarrow$ smoother.

- Bias-Variance Decomposition.

$$E[(f_{\theta(D)}(x') - y')^2] = (\bar{f}(x') - f(x'))^2 + E[(f_{\theta(D)}(x') - \bar{f}(x'))^2] + \sigma^2 \quad \text{where } \bar{f}(x') = E[f_{\theta(D)}(x')].$$

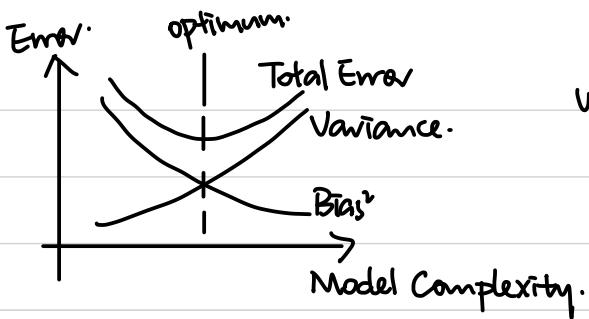
bias.

Variance

irreducible error

overfitting: learn different params for similar training set. \rightarrow High Variance

underfitting: learn similar params for different training set. \rightarrow High Bias

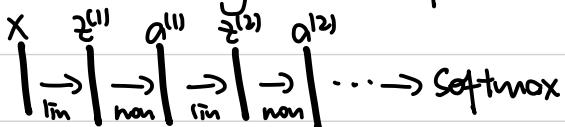


We don't fear overfitting.

Neural Network.

- Can represent nonlinear fns. Add nonlinearity = Activation fn.[↑]

Ex. $\tanh(z)$. Sigmoid $\sigma(z) = \frac{1}{1 + e^{-z}}$, $\text{ReLU}(z) = \max(0, z)$.



$w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}, \dots, w^{(L)}, b^{(L)}, w^{\text{final}}, b^{\text{final}}$ L hidden layers.

$f(x) = \text{softmax}(A^{\text{final}}(\sigma(A^{(1)}(\dots \sigma(A^{(L)}x)) \dots)))$

$$A^i(v) = W^i v + b^i$$

- Calculating gradients $\frac{\partial f}{\partial x_i} \approx \lim_{\epsilon \rightarrow 0} \frac{f(x + \epsilon e_i) - f(x - \epsilon e_i)}{2\epsilon}$

backpropagation (matrix-vector products).

Example: 2-layer

$$z^{(1)} = w^{(1)}x_i + b^{(1)} \quad a^{(1)} = \sigma(z^{(1)}) \quad z^{(2)} = w^{(2)}a^{(1)} + b^{(2)} \quad a^{(2)} = \sigma(z^{(2)}) \quad z = w^{\text{final}}a^{(2)} + b^{\text{final}}$$

$$\text{poly}_i(x_i) = \frac{\exp z}{\sum \exp z} \quad \log \text{poly}_i(x_i) = z_i - \log \sum \exp z \quad \ell(\theta; x_i, y_i) = \log \sum \exp z - z_i$$

$$\nabla z = \frac{\exp z}{\sum \exp z} - \delta_{yi} \quad \nabla a^{(2)} \ell = \frac{d^2}{da^{(2)}} \nabla z \ell = w^{\text{final}} \nabla z \ell.$$

$$\nabla w^{\text{final}} \ell = \frac{d^2}{dw^{\text{final}}} \nabla z \ell = (\nabla z \ell) A^{(2)\top} \quad \nabla b^{\text{final}} \ell = \frac{d^2}{db^{\text{final}}} \nabla z \ell = \nabla z \ell.$$

$$\nabla z^{(2)} \ell = \frac{d^2}{dz^{(2)}} \nabla a^{(2)} \ell = [\sigma'(z^{(1)}) \dots \sigma'(z^{(L)})] \nabla a^{(2)} \ell. \quad \nabla a^{(1)} \ell = \frac{d^2}{da^{(1)}} \nabla z^{(2)} \ell = W^{(2)\top} \nabla z^{(2)} \ell.$$

$$\nabla w^{(2)} \ell = \frac{d^2}{dw^{(2)}} \nabla z^{(2)} \ell = (\nabla z^{(2)} \ell) A^{(1)\top} \quad \nabla b^{(2)} \ell = \frac{d^2}{db^{(2)}} \nabla z^{(2)} \ell = \nabla z^{(2)} \ell.$$

$$\text{jacobian } J_{ij} = \frac{\partial f_i}{\partial x_j} \quad \text{Hessian } H_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j}$$

AutoDiff

> Jacobian-Vector Product JVP (forward)
build one Jacobian column at a time.

$$\frac{\partial f}{\partial x_1} = Jf(x) \overset{R^n}{\downarrow} e_1, \frac{\partial f}{\partial x_2} = Jf(x) e_2, \dots, \frac{\partial f}{\partial x_n} = Jf(x) e_n.$$

$$(a \rightarrow b) \rightarrow (a, T_a) \rightarrow (b, T_b)$$

$O(1)$ times the FLOPs $O(n(m m_3 + m_3 m_2 + m_2 m_1 + m_1 n))$ $\hookrightarrow X$ memory.

> Vector-Jacobian Product VJP (backward)

build one Jacobian row at a time.

$$\nabla f_1(x) = e_1^T Jf(x), \nabla f_2(x) = e_2^T Jf(x), \dots, \nabla f_m(x) = e_m^T Jf(x)$$

$$(a \rightarrow b) \rightarrow a \rightarrow (b, CT b \rightarrow o CT a)$$

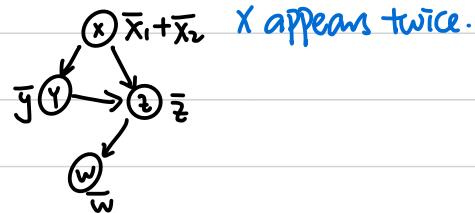
$O(1)$ times the FLOPs $O(m(m m_3 + m_3 m_2 + m_2 m_1 + m_1 n))$ $O(\text{depth})$ memory.



def $f(x)$:

$$y = \sin(x) \quad z = x * y \quad w = \cos(z)$$

return w .



Optimization

we don't really care about reaching the optimum.

- critical point \rightarrow global optimum / local optimum / saddle point / everywhere.

Stochastic Gradient Descent (SGD).

For each epoch, shuffle the data, construct mini batches B from consecutive data points,

$$\text{compute } \nabla_{\theta} \frac{1}{B} \sum_{i=1}^B p(\theta; x_i, y_i)$$

learning rate: high \rightarrow may cause oscillation / divergence. Low \rightarrow may stop learning too early.

(learning curve: epoch vs loss)

Learning rate schedule: Linear decay $\Delta t = \text{initial} \cdot (1 - \frac{i}{\text{max-steps}})$

Cosine decay $\Delta t = \text{initial} \cdot 0.5 \cdot [1 + \cos(\pi i + \frac{\pi}{\text{max-steps}})]$

Linear warming for large initial.

$\Delta \text{initial} = 0.001$ (need tuning) $k \Delta \text{initial} \rightarrow k B$.

- Momentum.

"remember" the steps before

$$\theta \leftarrow \theta - \alpha g. \text{ before: } g \leftarrow \nabla_{\theta} \frac{1}{N} \sum_{i=1}^N \ell(\theta; x_i, y_i) \text{ now: } g \leftarrow \theta_0 \frac{1}{N} \sum_{i=1}^N \ell(\theta; x_i, y_i) + \mu g.$$

exponential moving average

- Nesterov's accelerated gradient.

$$\theta \leftarrow \theta - \alpha g. \quad g \leftarrow \nabla_{\theta} \frac{1}{N} \sum_{i=1}^N \ell(\theta - \alpha g; x_i, y_i) + \mu g.$$

- Adam

(combine Momentum and second moment adjustment. Normalize the magnitude).

$$\theta \leftarrow \theta - \alpha g \quad \text{Momentum: } m \leftarrow (1 - \beta_1) \nabla_{\theta} l + \beta_1 m$$

$$E(x^t) \leftarrow \text{Second moment estimate: } v \leftarrow (1 - \beta_2) (\nabla_{\theta} l)^2 + \beta_2 v.$$

$$\text{detail-bias correction: } \hat{m} = m / (1 - \beta_1^t) \quad \hat{v} = v / (1 - \beta_2^t)$$

$$g = \hat{m} / (\sqrt{\hat{v}} + \epsilon)$$

- Adaptive learning rate (rescale separately for each dimension).

$$\theta_k^{t+1} = \theta_k^t - \frac{\alpha}{\sqrt{s_k^t + \epsilon}} \nabla_{\theta_k} L(\theta^t)$$

$$\text{RMSProp } s_k^t = \beta s_k^{t-1} + (1 - \beta) |\nabla_{\theta_k} L(\theta^t)|^2 \quad \text{moving average.}$$

$$\text{Adagrad } s_k^t = s_k^{t-1} + \beta (\nabla_{\theta_k} L(\theta^t))^2 \quad \text{Sum.} \rightarrow s^t \text{ monotonically increases.}$$

- Weight decay $\theta \leftarrow (1 - \lambda) \theta$.

for SGD, weight decay = L_2 -regularization.

for Adam, either weight decay (AdamW) or L_2 -regularization.

building-blocks

Normalization

- Input/Standardization. Preprocessing step / circularize the loss landscape. *Conti. only*

$$\text{for each dimension } d, \quad M_d = \frac{1}{N} \sum_{i=1}^N X_d \quad \sigma_d = \sqrt{\frac{1}{N} \sum_{i=1}^N (X_d - M_d)^2}$$

- normalizing activations.

Batch Normalization BN \rightarrow normalizing $\tilde{x}^{(l)}$ or $a^{(l)}$ by Batches. * large B

= add a BN layer before/after activation.

after normalization, each dimension * γ_d + β_d . scale shift.

< train mode

test mode. \rightarrow use statistics computed in train mode on one test sample.

* aggregated.

⑤ Leaky ReLU: $\text{ReLU}(x) = \begin{cases} x & x \geq 0 \\ \alpha x & x < 0 \end{cases}$ α is small

Convex

Exponential linear unit ⑥ $\text{ELU}(x) = \begin{cases} \alpha(x - 1) & x \leq 0 \\ x & x > 0 \end{cases}$

Convex

enables higher $\alpha \rightarrow$ faster training; training stability ↑; requires large B .

Layer Normalization: LN.

operates on single data points / by feature dimension, shared across dimension.

also learnable scales and shifts

- nonlinearity.

① Convex

$$\text{ReLU}(v) = \max\{0, v\}, \nabla_v \text{ReLU}(v) = \text{diag}(1_{[v>0]}) \leq 1 \quad \text{no gradient update for negative values.}$$

$$\text{② Sigmoid}(v) = \frac{\exp(v)}{1 + \exp(v)} = \frac{\exp(v)}{\exp(v) + 1} \quad \nabla_v \text{Sigmoid}(v) = \text{diag}[\text{Sigmoid}(v) \odot (1 - \text{Sigmoid}(v))] \leq \frac{1}{4}$$

Very small gradients for large values.

$$\text{③ Gaussian Error Linear Units: } \text{GELU}(v) = v \odot \tilde{\Phi}(v) \quad \tilde{\Phi}(v) \text{ cdf of } N(0, 1)$$

- skip connections = residual connections.

$$a^{(l)} = \sigma(z^{(l)}) + a^{(l-1)} \quad \text{Instead of } a^{(l)} = \sigma(z^{(l)}) \quad \text{better loss landscape.}$$

- weight Initialization.

$$x_j \sim N(0, 1) \quad E[z^2] = \sum_j E[(w^{(l)}_j)^2] E[x_j^2] = d \sigma_w^2$$

for every linear layer, $w_{ij} \sim N(0, \sigma_w)$ $\sigma_w^2 = \frac{1}{d} \rightarrow d \in \mathbb{N}$ $V_{\text{var}}(x) = E[x^2] - (E[x])^2$

- dropout: *inductive bias

training only: each iteration, randomly zero out p of w_{ij} .

- ensembles

Variance ↓, overfitting ↓, uncertainty quantification.

in practice: all, majority.

Convolutional Networks. * linear operator.

- filter $[k, k, 0, 1]$; bias - [0]

$$\frac{H+2P-K}{S} + 1 \quad \text{floor division}$$

input $[I, H, W]$, output $[O, H', W']$ $[H', W'] = \frac{[I, H, W] + 2 \times \text{pad} - (k-1)}{\text{stride}}$

- Convolution → BN | LN → ReLU

$[N, C, H, W]$ on N, H, W on C, H, W

- Pooling ↗ Max Pooling - non-linear

Average Pooling - linear

- Forward: $a^{(l)}([H, W, C]) \rightarrow z^{(l+1)}([H', W', C'])$ using filter $w^{(l+1)}([k, k, C', C])$

$$z^{(l+1)}[i, j, k] = \sum_{a=0}^{k-1} \sum_{b=0}^{k-1} \sum_{c=0}^{C-1} w^{(l+1)}[a, b, k, c] a^l[i+a, j+b, c]$$

$$\text{Matrix-Vector Notation: } z^{(l+1)}[i, j] = \sum_{a=0}^{k-1} \sum_{b=0}^{k-1} w^{(l+1)}[a, b] a^l[i+a, j+b]$$

$$\downarrow \quad \downarrow \quad \downarrow \\ c' \quad c' \times c \quad c$$

$$\begin{aligned} \text{- Backward: } & \frac{\partial z^{l+1}[i, j]}{\partial a^l[i+a, j+b]} = w^{l+1}[a, b] \quad \frac{\partial l}{\partial a^l[i, j]} = \sum_{a=0}^{k-1} \sum_{b=0}^{k-1} w^{l+1}[a, b] \frac{\partial l}{\partial z^{l+1}[i-a, j-b]} \\ & \frac{\partial z^{l+1}[i, j]}{\partial w^{l+1}[a, b]} = a^l[i+a, j+b] \quad \frac{\partial l}{\partial w^{l+1}[a, b]} = \sum_{i=0}^{H-1} \sum_{j=0}^{W-1} a^l[i+a, j+b] \frac{\partial l}{\partial z^{l+1}[i, j]} \end{aligned}$$

Computer Vision: e.g. flip, crop inductive bias
data augmentations invariances

the 3R's of Vision: Recognition, Reconstruction, Reorganization.

↓ attach semantics labels to objects ↓ obtain 3D info ↓ partition image based on semantics info.

Image classification

Object Localization: determine the bounding box.

Object Detection: measured by mean-average-Precision mAP; several objects; R-CNNs.

Semantics Segmentation: Label every pixel in the image; CNN, UNet.

Few challenges facing AI today: the limitations of supervision; static dataset AI VS embodied AI;

perception disconnected from action; the semantic gap.

the six lessons: be multi-modal; be incremental; be physical; explore; be social; use language.

*: empirical / training loss

Style transfer:

image generation: freeze parameters; update image.

input random image \leftrightarrow CNN \leftrightarrow V \leftrightarrow f(V).

output of intermediate conv layer.
(channel, height * width) 2D.

$$\text{Not exactly the same. } L_{SV} = \frac{1}{2} \sum_{ij} (V_{ij} - G_{ij})^2.$$

$$G_{ij} = \sum_k V_{ik} V_{jk}$$

$$L_{SV} \propto \sum_{ij} (G_{ij} - G_{ij}^*)^2$$

$$\ell = 2\beta c + \beta p_s \quad 2\beta = 10^{-3}.$$

Style loss operates on multiple intermediate activations of CNN with different weights. (Not f(c)).

Recurrent Networks.

Inputs: Sequential; Varying length. Outputs: anything / nothing.

$$\text{Weight sharing. } a^{(t+1)} = \sigma(W[a^{(t)}; X^{(t)}] + b). \quad * \begin{cases} h_t = \tanh((W_h h_{t-1} + W_x x_t) + B_h) \\ y_t = W_y h_t + B_y. \end{cases}$$

$a^{(0)}$ - independent from input.

Sequence input, Single output / Sequence input, sequence output. < Autoregressive

gradient < explode - clip.

gradient < vanish. - skip connection; LSTM. $\frac{de}{dw} = \frac{de}{da^t} \frac{da^t}{dw} + \frac{de}{da^{t-1}} \frac{da^{t-1}}{dw} + \dots + \frac{de}{da^1} \frac{da^1}{dw}$

large/small t. $\frac{de}{dw}$.

bidirectional models

LSTM: $f_t = \sigma(X_t U_f + h_{t-1} W_f)$ forget Gate. \leftarrow remember

$i_t = \sigma(X_t U_i + h_{t-1} W_i)$ input/update Gate.

$o_t = \sigma(X_t U_o + h_{t-1} W_o)$ output Gate.

$$\tilde{c}_t = \tanh(X_t U_g + h_{t-1} W_g).$$

memory. ↓

$$C_t = f_t C_{t-1} + i_t \circ \tilde{c}_t$$

output.

$$h_t = \tanh(C_t) \circ o_t.$$

$$\frac{\partial J}{\partial w_g} = \sum_{t=1}^T \frac{\partial J}{\partial h_t} \frac{\partial h_t}{\partial w_g} = \sum_{t=1}^T \frac{\partial J}{\partial h_t} \left(\frac{\partial h_t}{\partial C_t} \left(\frac{\partial C_t}{\partial w_g} + \frac{\partial C_t}{\partial C_{t-1}} \frac{\partial C_{t-1}}{\partial w_g} \right) + \frac{\partial h_t}{\partial C_{t-1}} \frac{\partial C_{t-1}}{\partial w_g} \right)$$

* (-ft)

1-1: classification, semantic segmentation.

1-n: image captioning

n-n: language translation.

Transformer.

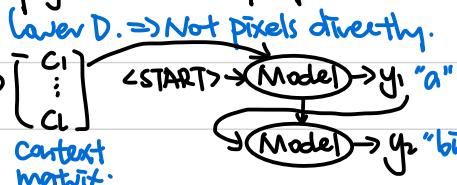
Sequential output ("A bird flying over a body of water.") one step at a time.

Attention: (Luong)

X → top layers of Pretrained Conv Net

(Layer D.) → Not pixels directly.

Context matrix.



$$q_2 \text{ want } k \text{ that is most similar to } q_2$$

$$k_1 \rightarrow e_{1,2} = k_1^T q_2 \rightarrow d_{1,2}$$

$$\vdots$$

$$k_L \rightarrow e_{L,2} = k_L^T q_2 \rightarrow d_{L,2}$$

$$\text{Softmax}_6 : = \sum_l a_{l,2} v_l.$$

$$\text{Similarity Score.} \downarrow$$

$$\text{Normalized.}$$

↑
Not y .

key-value-query: $q_t = q(\text{model info})$. e.g. want subject.

(all learned). $k_t = k(C_t)$ e.g. $(k_1, v_1) = (\text{subject}, \text{bird})$. $(k_2, v_2) = (\text{background}, \text{water})$.

$$V_t = V(C_t)$$

*linear.

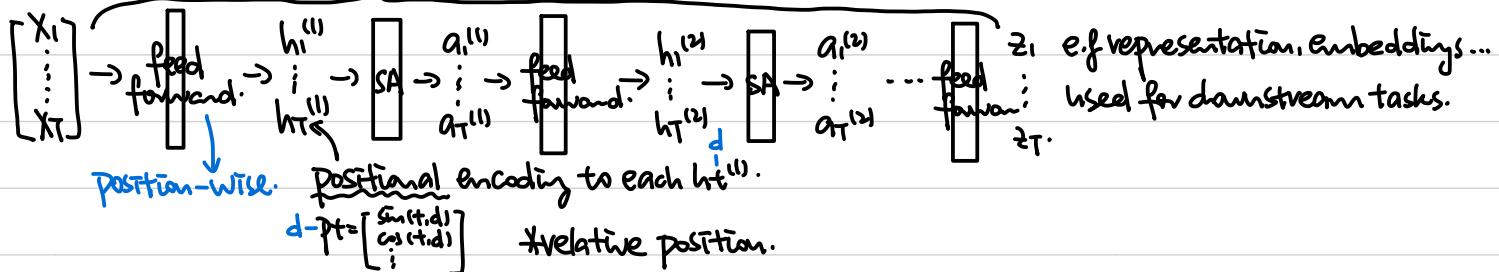
Self-Attention as a Layer: (handle sequential input). \Rightarrow processing whole sequence at once.

$$\begin{aligned} d-q_t &= q(X_t) & X_t \rightarrow q_t \\ d-k_t &= k(X_t) & k_t \rightarrow e_{t,2} = k_t^T q_t & \alpha_{t,2} \\ d-v_t &= V(X_t) & k_t \rightarrow e_{t,2} = k_t^T q_t & \alpha_{t,2} \\ & & & \alpha_t = \sum_i \alpha_{t,i} v_i \end{aligned}$$

$$\begin{array}{c} X_1 \\ \vdots \\ X_T \end{array} \xrightarrow{\text{SA}} \begin{array}{c} \alpha_1 \\ \vdots \\ \alpha_T \end{array}$$

Scaled dot product attention: $\frac{e_{t,2}}{\sqrt{d}}$

Transformer Encoder:



Multi-Head Attention: if want to capture multiple info at a time.

Independent heads \rightarrow concatenate to get a_2 .

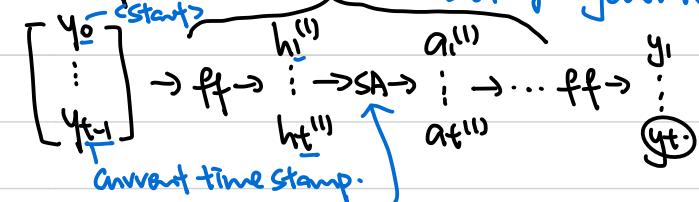
if $D=512$ for one head, scale down $D=64$ for q, k, v for 8 heads.

Transformers: optimization is hard: B+, LN.

LN, dropout, skip connection.

forward: $X + \text{self-droppout}(\text{Sublayer}(\text{Self-Attn}(X)))$.

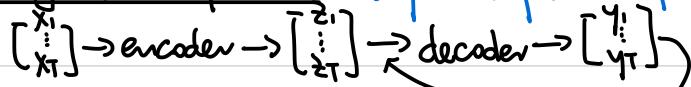
Transformer Decoder: training - fed truth, testing - generate $y_1 \dots y_T \rightarrow$ fed back.



Masked Attention: training. prevent from seeing the future.

e.g. $e_{T,2} = -\infty \Rightarrow \alpha_{T,2} = 0$. by MA, we can train T tokens simultaneously.

① Original transformer seq-to-seq model for translation.



② BERT: encoder only.

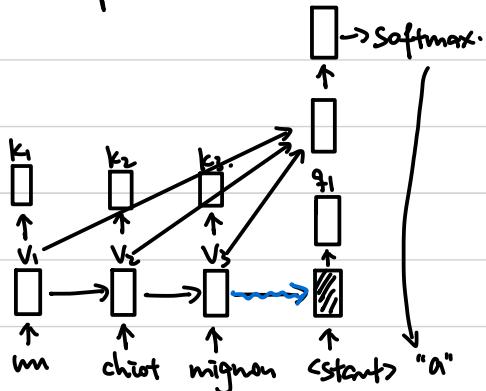
GLUE benchmark. (GLUE tasks).

- Pre-training: data + ttt; <CLS> <Sentence 1> <SEP> <Sentence 2>
 fine tuning: data - ; specific task.
- word-level representation.
 randomly mask 15% tokens → predict them.
 Predict the order of pair sentence. (binary classifier).
 sentence-level representation.
- ③ Vision transformers (ViTs): encoder only.
 ↗ no overlap, 1bx1b. ↗ no need to generalize.
 image → patches as sep. + positional encoding (learned)
- Swin: shifted window to perform SA.
 unsupervised ViT: train 2 ViTs to be similar.
 masked autoencoder: mask 75% → encoder → 25% + 75% empty → "decoder"
 MAE.
- ④ Reinforcement Learning. GPT-3 decoder only.
 parameters + ttt; good at answering questions, summary.

Seq2Seq. encoder-decoder Machine translation.

basic RNN (LSTM) seq2seq. > bottleneck problem.

Multi-layer RNNs.



Cross Attention.

key,value from encoder. query from decoder.

IS

trained on Ct; text input, text output.

learns the prompt along with the text.

Natural Language Processing.

M - size of vocab. T - max length of seq.

greedy encoding. - extend single hypothesis. $O(MT)$.

Beam Search. - extend multiple hypothesis. Prune.

every t, consider all possibilities, keep top k (largest p) only.
e.g.t. low p to all words.

prefer short ; length normalization: $\frac{S(e)}{m}$. OR add bonus for extra word.

Label Smoothing: 90%, 10%.

$$S(e) = \sum_{i=1}^m \log p(e|e_i, f_i)$$

Byte-Pair Encoder. Instead of <UNK>.

Back Translation.

Unsupervised pretraining. (get embedding).

trained on some task, use hidden state as embedding.

ELMo: bidirectional LSTM, concatenate hidden states.

GPT: transformer-based forward language model (past context only); text generation.

BERT: (bidirectional models are inherently unidirectional.) \Rightarrow Single transformer encoder
incorporate bi context. ; No masked attention.

T5: BERT style mask, but encoder-decoder ; reform all to unified text-to-text format.
more general and flexible ; multi-task learning.

Transformer over LSTM: model long-term dependencies and parallelization.

Distribution shift.

when test set has a different distribution than training set.

Risk.

Empirical Risk (ERM). training only. generalization bounds. The big lie of ML.

Benchmark: ① ImageNet challenge test sets.

measure robustness under different DS ; stress test ; same classes.

② Wilds - different domains.

task of premise-hypothesis, contradiction, neutrality, entailment.

③ in NLP, Adversarial natural language inference (ANLI) dataset.

Against DS: ① training larger models on larger datasets.

② Data Augmentation.

- Mixup: element-wise convex combinations of data points. from \neq images.
$$(x_1 + (1-\alpha)x_2, y_1 + (1-\alpha)y_2)$$
- AutoAugment: learns augmentation strategies.
- AugMix: mix random augmentation. Autoaugmentation.
- PixMix: mix in separate image datasets

③ Alternative / additional training objectives.

ViT: masked autoencoding objective.

Anomaly: Anomaly score, generative model. ($\log p(x) - \text{Anomaly}$). x

Detection: Simple baseline: Confidence = $\max_k P(\text{y}=k|x)$ Anomaly score = $-\max_k P(\text{y}=k|x)$.
(out-of-distribution). cannot detect adversarial examples.

Benchmark: train a model on one dataset, and treat any other datasets as anomaly.

e.g. ImageNet-22k, species.

binary classification problem. TP TN FP FN.

Model calibration: - Compare confidence against accuracy. $E[\mathbb{P}(Y=\hat{Y}|\hat{p}=p) - p]$.

e.g. 60% confidence raining \Rightarrow rain 60% of time. $P(Y=\hat{Y}|\hat{p}=p)=p \quad \forall p \in [0,1]$

① ensemble of models.

② (Single model) Add a softmax temperature T.

$\hat{p}(y=\hat{y}|x) = \frac{\exp(\hat{p}_i/T)}{\sum_{j=1}^k \exp(\hat{p}_j/T)}$ T-hyperparameter. only calibration.

minimize negative log-likelihood.

> Expected Calibration Error (ECE).

group predictions into M bins. $B_m = \{(y_i, \hat{p}_i) \text{ s.t. } \hat{p}_i \in (\frac{m-1}{M}, \frac{m}{M}]$.

$$\text{acc}(B_m) = \frac{1}{|B_m|} \sum_{i \in B_m} \mathbb{1}\{\hat{y}_i = y_i\} \quad \text{conf}(B_m) = \frac{1}{|B_m|} \sum_{i \in B_m} \hat{p}_i$$

$$\text{ECE} = \sum_{m=1}^M \frac{|B_m|}{n} |\text{acc}(B_m) - \text{conf}(B_m)|.$$

> Maximum Calibration Error (MCE).

$$\text{MCE} = \max_{m \in \{1, \dots, M\}} |\text{acc}(B_m) - \text{conf}(B_m)|.$$

Robustness.

Handle DS w/o lots of data, domain knowledge.

① Domain Adaptation.

Assumption: abundant data in Source Domain, few data from target domain.

Approaches: ① Importance Weighting. (Similar ones get higher weights).

$$E_{\text{target}}[\ell(\theta; x, y)] = E_{\text{source}} \left[\frac{P_{\text{target}}(x, y)}{P_{\text{source}}(x, y)} \ell(\theta; x, y) \right].$$

Unsupervised domain adaptation: No y from target domain

assume covariate shift: $P_{\text{target}}(y|x) = P_{\text{source}}(y|x)$.

② Invariant feature learning (Want features \leftarrow informative \Rightarrow predict y .
domain discriminators: \leftarrow invariant \rightarrow same distribution for S & T).

$x \rightarrow \text{features} \rightarrow \text{class label } y$.

\leftarrow gradient reversal layer \rightarrow domain label d .

Invariant Risk Minimization (IRM) optimal classifiers are same.

directly make feature distribution similar

Correlation Alignment (CORAL). - Add a loss term.

② Subpopulation Shift.

Assumption: Not only S & T. Some domains are underrepresented in training data, but at test time, we want fairness / weight / worst-case ("Robustness").

In distribution (group) robustness, Adversary can only change distribution across domains, not distribution of classes inside domain. $\min_{\theta} \max_{D \in \mathcal{D}} \sum_{d=1}^D P_d \mathbb{E}_d [\ell(\theta; x, y)]$.

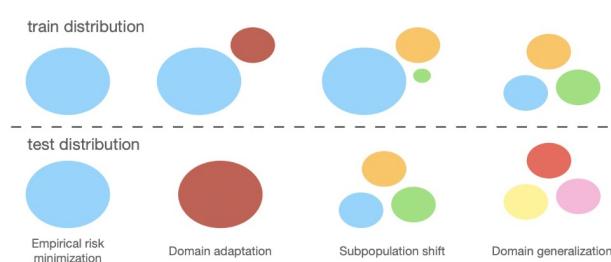
rebalancing the training data (upsampling / down sampling).

③ Domain Generalization. (One-shot domain Adaptation / Multi-Source domain Adaptation).

Assumption: > 2 domains; No domain imbalance; will be given new domain at the test time \Rightarrow want to generalize.

Summary: different assumptions on training, expected test shift.

"domains" \Rightarrow share.



Test-time-adaption: change prediction at test time. e.g. threshold, rotation.
 assume multiple test points available.
 Adaptive model: $x \xrightarrow{p_x} y$. distribution of where x comes from.

Not distribution shift.

Adversarial Robustness.

modern adversarial distortions: perceptible to human eyes but not perceptible to models.

ℓ_1 Norm: $\|V\|_1 = |V_1| + |V_2| + \dots + |V_d|$

$$\|V\|_2 = \sqrt{|V_1|^2 + \dots + |V_d|^2}$$

$\|V\|_\infty = \max\{|V_1|, |V_2|, \dots, |V_d|\}$. greatest change among dimensions.

Cumulative effect of many small changes made the adversary powerful. (change class of \hat{y}).

distortion budget ϵ s.t. $\|x_{adv} - x\|_p \leq \epsilon$. ℓ_p attack.

Adversary's goal: find distortion δ s.t. $x_{adv} = x + \underset{0 \leq \delta \leq \epsilon}{\operatorname{argmax}} \ell(\theta; x + \delta, y)$.

Fast Gradient Sign Method. (FGSM)

$$x_{FGSM} = x + \epsilon \operatorname{sign}(\nabla_x \ell(\theta; x, y)) \quad \|x_{FGSM} - x\|_\infty = \epsilon. \text{ Single step GD.}$$

Projected Gradient Descent (PGD)

Random initialize a perturbed image $\tilde{x} = x + n$, $n \sim \mathcal{U}(-\epsilon, \epsilon)$, initialize $\delta = 0$.

for $t=1, \dots, T$: $\delta \leftarrow \operatorname{clip}(\delta + \alpha \operatorname{sign}(\nabla_\delta \ell(\theta; \tilde{x} + \delta, y)), -\epsilon, \epsilon)$

$$x_{PGD} = \tilde{x} + \delta$$

Method 1: Adversarial Training. (AT). will reduce accuracy on nonadversarial examples.

sample batch $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}) \dots (x^{(B)}, y^{(B)})$ in training set.

Create $x_i^{(i)adv}$ (OR $x_i^{(i)PGD}$) from $x_i^{(i)}$ for all i .

optimize average training loss on them.

Method 2: More data.

Adversarial robustness scales slowly with dataset size.

✓ adversarial pretraining on larger dataset.

Method 3: Data Augmentation + adversarial training + parameter exponential moving average \Rightarrow CutMix. (MixUp + CutOut).

untargeted attack. - maximize the loss $\underset{\delta: \| \delta \|_p \leq \epsilon}{\operatorname{argmax}} \ell(\theta; x + \delta, y)$

targeted attack. - optimize examples to be misclassified as predetermined \hat{y} .
 $\underset{\delta: \| \delta \|_p \leq \epsilon}{\operatorname{argmin}} \ell(\theta; x + \delta, \hat{y})$.

"arms race" that defenders lose; within weeks.

transferability: models with different architecture. Real-world instantiation noise, sensor noise.

Smooth GELU instead of ReLU.

robustness to unforeseen attacks.

Deep generative Models.

① Generation: synthesize new data points.
 (Conditional generation).

② Density Modeling.

③ Representation Learning. latent representation z .

Generative adversarial networks (GANs) ①③ no need to compute $P(x)$.

generator G - synthesize data

vs. discriminator D - distinguish real data from synthesized data. (throw away afterwards).

$$\min_{G} \max_{D} \mathbb{E}_{x \sim P_{\text{data}}} [\log D(x)] + \mathbb{E}_{z \sim p(z)} [\log (1 - D(G(z)))] \quad x = G(z) \sim P_G(x)$$

optimal G^* under Bayes optimal D^* recovers the data distribution.

$$D^*(x) = \frac{P(x)}{P(x) + P_G(x)} \Rightarrow \min_{G^*} \mathbb{E}_{x \sim P(x)} [\log (\frac{P(x)}{P(x) + P_G(x)})] + \mathbb{E}_{x \sim P_G(x)} [\log (\frac{P_G(x)}{P(x) + P_G(x)})].$$

$$\text{(let } q(x) = \frac{P(x) + P_G(x)}{2} \text{)} \quad \min_{G^*} \mathbb{E}_{x \sim P(x)} [\log P(x) - \log q(x)] + \mathbb{E}_{x \sim P_G(x)} [\log P_G(x) - \log q(x)] + \text{constant.}$$

JSD divergence between $P(x)$ and $P_G(x)$. \Rightarrow generator distribution to close data distribution.

better architecture, modified objectives, bags of tricks, scaling up.

Conditional generation: class label OR "label-like" input (infoGAN) to both G and D

cycleGAN: $2G, 2D$. turns image from one domain to another, reverse.

alternatively taking GD on G and D . hard to evaluate (pretrained classifier)

In practice, care about whether D provides a useful signal to G . \Rightarrow restrict D 's expressivity.

↑ Smooth: ① Lipschitz constraint. ② real-valued D. ③ instance noise

$$W(P, P_G) = \sup_{\|f\|_1 \leq 1} [E_{x \sim P}(f(x)) - E_{x \sim P_G(x)}(f(x))]. \quad f: 1\text{-Lipschitz}$$

① weight clipping. ($|$ each entry of weight $| < \epsilon$). k -Lipschitz.

② gradient penalty. $+ \lambda (\|\nabla_x f(x)\|_2 - 1)^2$ \sim 1-Lipschitz.

③ spectral Normalization. $\text{ReLU} \rightarrow$ 1-Lipschitz. $k = \sup_{\|x\|_2=1} \|Wx\|_2 = \sigma(W)$ $W_i \leftarrow \frac{W_i}{\sigma(W)}$ largest singular value of W

Autoregressive Models ①, ②.

$P(x) = \prod_{i=1}^d P(x_i | x_{\leq i})$ masked convolution / masked SA. pixels \rightarrow mixture of logistic distributions
Very slow.



Latent Variable Models

known: $P(z)$, with Gaussian $P(x|z)$. Intractable: $P(x) = \int p(z) p_\theta(x|z) dz$. $p(z|x)$.

model $P(x|z)$ by $q_\phi(z|x)$. \Rightarrow minimize KL divergence.

$$\begin{aligned} \text{DKL}(q_\phi||P_\theta) &= E_{q_\phi}[\log q_\phi(z|x) - \log p_\theta(z|x)]. = E_{q_\phi}[\log q_\phi(z|x) - \log p_\theta(x|z)] + \log p_\theta(x). \\ &= E_{q_\phi}[-\log p_\theta(x|z) + \log q_\phi(z|x) - \log p(z) + \log p_\theta(x)] = -[E[\log p_\theta(x|z)]] - \text{DKL}(q_\phi||P_\theta) + \log p_\theta(x) \geq 0. \\ \text{objective: } \arg \max_{\theta} \sum_i \log p_\theta(x_i) &= \arg \max_{\theta} \sum_i \log \left(\sum_z p(z) p_\theta(x_i|z) \right). \end{aligned}$$

reconstruction + regularization

ELBO · lower bound $\nearrow \phi$

$$\text{by } (\log E[x] \geq E[\log(x)]), \geq \sum_{i=1}^n E_{q_\phi(z|x_i)} [\log p(z) - \log q_\phi(z|x_i) + \log p_\theta(x_i|z)].$$

Re-parametrization Trick: $q_\theta(x) \sim N(\mu, \Sigma)$ $E_q(x^*) = E_p[(\mu + \Sigma Z^2)]$. $\Sigma \in N(0, I)$. $\mu \sim N(0, I)$.

> Variational autoencoders (VAEs). ① ② ③ \rightarrow ELBO.

observation model $P(x|z)$, recognition model $q_\phi(z|x)$ $\xrightarrow{\text{encoder-decoder}}$ minimize ELBO

Variational inference (VI) framework. keep both.

blurred. \leftarrow conditional independence assumption.

> deep autoencoders. (e.g. denoising, bottleneck).

> masked autoencoders. ③, little ①.

Diffusion Models

fix forward process (x to z): add Gaussian noise, little by little \Rightarrow z (pure noise).

reverse process: (z to x).

maximizing ELBO.

e.g. DALL-E 2.

deep generative models
Self-supervised. \rightarrow deep unsupervised learning.
 \hookrightarrow target representation learning.

Pretend a part of input is unknown and predict it. e.g. BERT; MAE

① leverage known structure of images. ② contrastive learning. ③ data augmentation.
leverage spatial context; predicting rotations.

Contrastive Learning: learned representation should be close to "similar". $\ell = -\log \frac{\exp(\mathbf{z}^T \mathbf{z}_j + \mathbf{c})}{\sum_{k=1}^K \exp(\mathbf{z}^T \mathbf{z}_k + \mathbf{c})}$ similar.

Contrastive predictive Coding (CPC): similar \rightarrow comes from same input; extracting.

Similar: diff. augmentation of the same image. SimCLR: on $g(\cdot)$ not $f(\cdot)$. additional transformation representation.

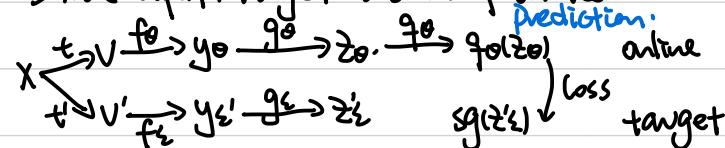
Momentum Contrast (MoCo): queue of negatives (add representation from the latest mini batch.).

EMA; encourages similarity across minibatch.

technical improvements \Rightarrow SimCLR: $\ell_{i,j} = -\log \frac{\exp(\mathbf{z}_i^T \mathbf{z}_j + \mathbf{c})}{\sum_{k=1}^N \exp(\mathbf{z}_i^T \mathbf{z}_k + \mathbf{c})}$

/ Bootstrap your own latent (BYOL): online vs. target network; given two diff. augmentation of degenerate solution.

Same input; target is EMA of online.



nothing.

Self-distillation with no labels (DINO): student \rightarrow teacher's prob; "distillation": teacher knows sth.

Multimodal contrastive learning (CLIP): text encoder & image decoder; similar: (text, image).

Massive Models (transformer decoder (language models)).

Scaling laws: larger models require fewer samples. ; 8x model size, 5x dataset size.
(neural language models). train large models, stop short after convergence.

GPT-3: few-shot learning. #examples.

Gopher: 280B. 2T \rightarrow 300B.

Chinchilla: varying hyperparameters; X model = X data.; Massive Multitask Language Understanding (MMLU)

Megatron-Turing NLG: 530B. 339B \rightarrow 270B.

PaLM: 540B. 780B \rightarrow 780B.

chain-of-thought prompting. $\downarrow^{<10}$ self-consistency. (sampling + pick answer).
in-context learning.
process to achieve answer.

Fine Tuning: update via gradient based optimization with small datasets.

Specializing Code: Codex, AlphaCode.

Scaling up model / dataset size; change architecture to encoder-decoder; fine tune on competition code; sample/filter sols.

Meta-Learning.

Assumption: meta-training task & meta-test task drawn i.i.d. from same distribution.

Input: D_{tr} , X_{ts}^{ts} Output: y_{ts}^{ts} .

Multi-task learning. - Solve T_1, \dots, T_n at once.

Transfer learning. - Solve T_b after solving T_a , transfer knowledge.] impractical to prior tasks.

meta-learning - given data T_1, \dots, T_n , quickly solve T_{test} .

/ query set.
Support set

K-shot learning, N-way classification.

① Black-Box Meta-learning.

Parametrize learners as NN. + expressive. - optimization not consistent.

② Optimization-based Meta-learning.

embed optimization. - 2nd optimization. Compute. + consistent, positive inductive bias.

③ Non-parametric.

+ expressive. Compute

→ Shot: independently, aggregate by $C_n = \frac{1}{k} \sum_{(x_i, y_i) \in D_{tr}} \mathbb{1}(y_i = n) f(x_i)$. $Poly = u(x) = \frac{\exp(-d(f(x), C_n))}{\sum_i \exp(-d(f(x), C_i))}$

meta-training: sample task T_i .

binary classification: sample disjoint D_{tr}, D_{test} from D_i .

Compute $\phi_i \leftarrow f(D_i^{tr})$ $y_{ts}^{ts} = \sum_{x_k, y_k \in D_{tr}} f(x_k, y_k)$

update ϕ_i by $DLL(\phi_i, D_i^{test})$. $DLL(y_{ts}^{ts}, y_{ts}^{ts})$.

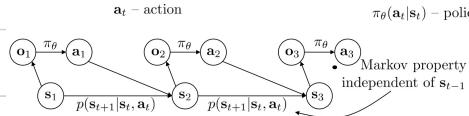
meta-testing: given D_j^{tr} from new task T_j .

N-way classification: Compute $\phi_j \leftarrow f(D_j^{tr})$.

Deep-Learning.

Moravec's Paradox.

s_t - state
 o_t - observation
 a_t - action
 $\pi_\theta(a_t|o_t)$ - policy
 $\pi_\theta(a_t|s_t)$ - policy (fully observed)



transition operator / probability.

goal:

$$p_\theta(s_1, a_1, \dots, s_T, a_T) = p(s_1) \prod_{t=1}^T \pi_\theta(a_t | s_t) p(s_{t+1} | s_t, a_t)$$

$$p_\theta(\tau)$$

$$\theta^* = \arg \max_{\theta} E_{\tau \sim p_\theta(\tau)} \left[\sum_t r(s_t, a_t) \right]$$

$$\text{RL objective: } \max_{\pi} \sum_{t=1}^T E_{s_t, a_t \sim \pi}[r(s_t, a_t)]$$

$$\text{Q-function: } Q^\pi(s_t, a_t) = \sum_{t'=t}^T E_{s_{t'}, a_{t'} \sim \pi}[r(s_{t'}, a_{t'}) | s_t, a_t]$$

$$\pi(a|s) = 1 \text{ if } a = \arg \max_a Q^\pi(s, a)$$

$$Q^*(s, a) = r(s, a) + \max_{a'} Q^*(s', a')$$

enforce this equation at all states!

$$\text{minimize } \sum_i (Q(s_i, a_i) - [r(s_i, a_i) + \max_{a'_i} Q(s'_i, a'_i)])^2$$

$$\text{minimize } \sum_i (Q(s_i, a_i) - y_i)^2$$

Behavioral cloning.

Dagger.

1. train $\pi_\theta(a_t|o_t)$ from human data $\mathcal{D} = \{o_1, a_1, \dots, o_N, a_N\}$
 2. run $\pi_\theta(a_t|o_t)$ to get dataset $\mathcal{D}_\pi = \{o_1, \dots, o_M\}$
 3. Ask human to label \mathcal{D}_π with actions a_t
 4. Aggregate: $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_\pi$

Q-learning:

$$Q(s, a) \leftarrow r(s, a) + \max_{a'} Q(s', a')$$

off-policy Q-learning:

1. collect dataset $\{(s_i, a_i, s'_i, r_i)\}$ using some policy, add it to \mathcal{B}
 2. sample a batch (s_i, a_i, s'_i, r_i) from \mathcal{B}
 3. minimize $\sum_i (Q(s_i, a_i) - [r(s_i, a_i) + \max_{a'_i} Q(s'_i, a'_i)])^2$