

linear models: loss: $\hat{y} \uparrow \text{cost } \varepsilon = \frac{1}{N} \sum_i (y_i - \hat{y}_i)^2$

① Linear regression: $y = Wx + b$ ($x: dx_1$) $\hat{y} = XW + b$ ($x: Nx_1$) **optimization**: $L(W) = \frac{1}{N} \|XW - \hat{y}\|^2$ $\nabla_W L = \frac{1}{N} X^T (XW - \hat{y})$ $W^* = (X^T X)^{-1} X^T \hat{y}$ **local minima**: nn is not convex (weight space symmetry)

Each GD cost O(ND), $(\varepsilon)^d \text{ cost } O(D^3)$ \Rightarrow random start, local minima are generally fine

② Logistic regression: $Z = W^T X$ $y = \sigma(Z) = \frac{1}{1+e^{-Z}}$ $L_{CE} = -t \log y - (1-t) \log(1-y)$ as two h have same incoming & outgoing weights squared gradient. An exponential moving average of the squared gradient $S_j \leftarrow (1-\gamma)S_j + \gamma \frac{\partial J}{\partial \theta_j}^2$ $\theta_j \leftarrow \theta_j - \alpha \frac{\partial J}{\partial \theta_j}$

③ Softmax regression: $Z = W^T X$ $\hat{y} = \text{softmax}(Z)$ $y_k = \frac{e^{Z_k}}{\sum e^{Z_m}}$, $L = -\hat{y}^T \log(y)$ \rightarrow always same gradient \rightarrow always identical

multi-layer perceptrons: $X_1 \rightarrow D_1 \rightarrow \dots \rightarrow D_n \rightarrow Y$

Convex: $\nabla \text{convex} \rightarrow$ if $x, x_2 \in S$, $\lambda x_1 + (1-\lambda)x_2 \in S$

④ $f((1-\lambda)x_1 + \lambda x_2) \leq (1-\lambda)f(x_1) + \lambda f(x_2)$ $\nabla \text{lower cost better perf}$

⑤ nn: permutation symmetries

⑥ local min is global min (strict = only 1 global min)

Fully connected layer: $y_{ni} = \phi(Wx_i + b)$

ReLU: ∇_x non differentiable at $x=0$

soft ReLU: $\sigma(z) = \log(1+e^z)$ ∇_z comp expensive

Expressive Power: if all activation func are linear \rightarrow MN can only represent linear function, no more than linear expression

Multi-layer feed forward nn with non linear activation func \rightarrow universal function approximator = approximate any f $\exists \epsilon, \delta$

eg: binary input & output: each hidden unit is 1 situation

Back propagation: $\frac{\partial y}{\partial z} = \begin{pmatrix} \frac{\partial y_1}{\partial z_1} & \frac{\partial y_1}{\partial z_2} \\ \vdots & \vdots \\ \frac{\partial y_n}{\partial z_1} & \frac{\partial y_n}{\partial z_2} \end{pmatrix} = \frac{\partial y}{\partial z} \cdot \frac{\partial z}{\partial z_1} \cdot \frac{\partial z}{\partial z_2} \cdots \cdot \frac{\partial z}{\partial z_n} = \frac{\partial y}{\partial z} \cdot \frac{\partial z}{\partial x} = \frac{\partial y}{\partial x}$ quadratic time

$y = \exp(z) \quad \frac{\partial y}{\partial z} = \frac{\partial y}{\partial z} \text{ is diagonal} \quad \bar{z} = \bar{y} \circ \exp(z) \text{ WTSE}$

$H \rightarrow y, y = Wh + b \quad \bar{w} = \bar{y} \cdot h^T \quad \bar{h} = w^T \bar{y}$

$W_{ki} = \bar{y}_k h_i, \bar{h}_i = \sum_k \bar{y}_k W_{ki}$ two add-mul per weight

the BP is expensive as $2FP$

For multi-layer perceptron, the cost is linear in # of layers and quadratic in # of units ($\propto k^2$ weights)

Auto Diff: $\frac{\partial y}{\partial x_1} = \frac{\partial y}{\partial x_1} \frac{\partial x_1}{\partial m_1} = \frac{\partial y}{\partial m_1} \frac{\partial m_1}{\partial x_1} = \dots = \frac{\partial y}{\partial m_N} \frac{\partial m_N}{\partial x_1}$

input \rightarrow output defnvp / exp, lambda g, ans, x: g * ans

$h \rightarrow \frac{\partial h}{\partial x} \quad y = \frac{\partial y}{\partial h} \quad (\bar{y}) \quad (y) \quad (\bar{h} = y \bar{x})$

$w \rightarrow \frac{\partial w}{\partial x} \quad Z = \frac{\partial Z}{\partial x} = \frac{\partial Z}{\partial h} + \frac{\partial Z}{\partial t} \quad \# \text{ of connections} = \frac{\partial Z}{\partial h} + \frac{\partial Z}{\partial t}$

$b \rightarrow \frac{\partial b}{\partial x} \quad \# \text{ of weights} = \frac{\partial b}{\partial h}$

Distributed representation: language modeling \approx predict the next word

of prev words: deeper network w/o significantly more computation complex

Markov Assumption: prob of the next word depends on a fixed optimization and Generation

$P(W_t | W_1, W_2, \dots, W_{t-1}) = P(W_t | W_{t-1}, W_{t-2}, \dots, W_1)$ ∇ stochastic gradient descent

N-Gram model: localist representation \rightarrow not share btw N-grams. Training curve (loss vs # of iteration or epoch)

Neural language model: distributed representation: words are represented as vectors in a continuous space.

Pros: semantic similarity, generalization (unseen data), robustness, efficiency (mem, comp), scalability (data, task)

Loss: $-\frac{1}{N} \sum_i t \log y_i$ (t -one hot encoder 1 for correct word, 0 elsewhere)

vs localist: share info between entity, complex to interpret

capture more nuanced info, robust to damage, noise, scale index of word \rightarrow embedding \rightarrow hidden unit \rightarrow softmax

Global Vector Embedding (Glove): skip layer connection

occurrence matrix: frequency of 2 words appear nearby

$X = R \times R^T$, $x_{ij} = y_i^T \tilde{y}_j$, $f(x_{ij}) = \frac{1}{(x_{ij}/x_{max})^d}$, $0 \leq f(x_{ij}) \leq 1$, if $x_{ij} > x_{max}$

$J(R) = \sum_{i,j} f(x_{ij})(y_i^T \tilde{y}_j + b_i + b_j - \log x_{ij})$, \log diminish high count importance, reduce count range, emphasize relative diff, balance the influence of diff count, s.t. frequent, infrequent pairs contribute meaningfully. SGD struggles when the loss is steep vertically but flat horizontally, progress slowly along the flat direction, oscillates wildly along steep

Pros: global statistics, semantic rela, quick trained, pre train ava

Cons: static embedding, x-sign mem for store & comp, unseen w in train

Momentum: Add a fraction of the previous update to the current update $m_{k+1} \leftarrow \mu m_k + \alpha \frac{\partial J}{\partial \theta}$ smooth out the updates strategy to counter LR decay = $\exp T LR$ each update

LR damping factor: $\theta \leftarrow \theta - m_{k+1}$ \Rightarrow how quickly the contribution of prev gradient decay.

RMS prop (Root Mean Square propagation): Adjust the learning rate for each parameter j based on the moving avg S_j of the squared gradient $S_j \leftarrow (1-\gamma)S_j + \gamma \frac{\partial J}{\partial \theta_j}^2$ $\theta_j \leftarrow \theta_j - \alpha \frac{\partial J}{\partial \theta_j}$

Weight decay: Normalize weight to unit vector after weight decay. Adam with L2 regularization

language modelling using RNN: training with teacher forcing = use target output as ground truth as input for next time step. Generating use model output as inputs for next time step. Each word is a one hot vector, predicts a dist. ∇ pmap over the next word. CE loss. long distance depend.

Neural Machine translation challenge: word order, sentence length, sequence to sequence architecture, future words, full context

two distinct networks with diff weights: Encoder - Process entire input sentence, compress into a code vector. Decoder - generate translation word by word

Exploding and vanishing gradients: RNN suffer during back prop through time, gradients are multiplied by the same weight matrix repeatedly. This repeated multiplication causes the gradient to explode or vanish.

Data augmentation: Augment training set (eg transformation), not test

Reduce Parameters: Add linear bottleneck layer (eg Autoencoder, Dimb)

Regularization: $\Sigma_{reg} = \varepsilon + \lambda R$ $\nabla_w \Sigma_{reg} = -\alpha \frac{\partial \Sigma_{reg}}{\partial w}$ $L_1: R = \frac{1}{L} \sum_i |w_i|$

Early stopping: $\# \text{ epochs}$

Stochastic Regularization: Avg or combine $\#$ of model prediction

Size: # of units: store activations in memory for back prop train

of weights \rightarrow # of parameters \rightarrow amount of overfitting train & test

of connections \rightarrow complexity of backprop (3 add multiply per connect) \rightarrow amount of overfitting train & test

Fully connected convolution: weight decay, a set of filters (kernel), a set of feature map

Pooling layer: avg, max, L2 norm, weighted avg

reduce dimension (L of Para) \rightarrow L comp, L speed, overfit

support translation invariance \rightarrow higher layer cover larger region

receptive field T (region of input image affect output) \rightarrow scale the weights by $1/T$ at test time to approximate the expected activation

size: # of units: store activations in memory for back prop train

of weights \rightarrow # of parameters \rightarrow amount of overfitting train & test

of connections \rightarrow complexity of backprop (3 add multiply per connect) \rightarrow amount of overfitting train & test

Fully connected convolution: weight decay, a set of filters (kernel), a set of feature map

Pooling layer: avg, max, L2 norm, weighted avg

reduce dimension (L of Para) \rightarrow L comp, L speed, overfit

support translation invariance \rightarrow higher layer cover larger region

receptive field T (region of input image affect output) \rightarrow scale the weights by $1/T$ at test time to approximate the expected activation

size: # of units: store activations in memory for back prop train

of weights \rightarrow # of parameters \rightarrow amount of overfitting train & test

of connections \rightarrow complexity of backprop (3 add multiply per connect) \rightarrow amount of overfitting train & test

Fully connected convolution: weight decay, a set of filters (kernel), a set of feature map

Pooling layer: avg, max, L2 norm, weighted avg

reduce dimension (L of Para) \rightarrow L comp, L speed, overfit

support translation invariance \rightarrow higher layer cover larger region

receptive field T (region of input image affect output) \rightarrow scale the weights by $1/T$ at test time to approximate the expected activation

size: # of units: store activations in memory for back prop train

of weights \rightarrow # of parameters \rightarrow amount of overfitting train & test

of connections \rightarrow complexity of backprop (3 add multiply per connect) \rightarrow amount of overfitting train & test

Fully connected convolution: weight decay, a set of filters (kernel), a set of feature map

Pooling layer: avg, max, L2 norm, weighted avg

reduce dimension (L of Para) \rightarrow L comp, L speed, overfit

support translation invariance \rightarrow higher layer cover larger region

receptive field T (region of input image affect output) \rightarrow scale the weights by $1/T$ at test time to approximate the expected activation

size: # of units: store activations in memory for back prop train

of weights \rightarrow # of parameters \rightarrow amount of overfitting train & test

of connections \rightarrow complexity of backprop (3 add multiply per connect) \rightarrow amount of overfitting train & test

Fully connected convolution: weight decay, a set of filters (kernel), a set of feature map

Pooling layer: avg, max, L2 norm, weighted avg

reduce dimension (L of Para) \rightarrow L comp, L speed, overfit

support translation invariance \rightarrow higher layer cover larger region

receptive field T (region of input image affect output) \rightarrow scale the weights by $1/T$ at test time to approximate the expected activation

size: # of units: store activations in memory for back prop train

of weights \rightarrow # of parameters \rightarrow amount of overfitting train & test

of connections \rightarrow complexity of backprop (3 add multiply per connect) \rightarrow amount of overfitting train & test

Fully connected convolution: weight decay, a set of filters (kernel), a set of feature map

Pooling layer: avg, max, L2 norm, weighted avg

reduce dimension (L of Para) \rightarrow L comp, L speed, overfit

support translation invariance \rightarrow higher layer cover larger region

receptive field T (region of input image affect output) \rightarrow scale the weights by $1/T$ at test time to approximate the expected activation

size: # of units: store activations in memory for back prop train

of weights \rightarrow # of parameters \rightarrow amount of overfitting train & test

of connections \rightarrow complexity of backprop (3 add multiply per connect) \rightarrow amount of overfitting train & test

Fully connected convolution: weight decay, a set of filters (kernel), a set of feature map

Pooling layer: avg, max, L2 norm, weighted avg

reduce dimension (L of Para) \rightarrow L comp, L speed, overfit

support translation invariance \rightarrow higher layer cover larger region

receptive field T (region of input image affect output) \rightarrow scale the weights by $1/T$ at test time to approximate the expected activation

size: # of units: store activations in memory for back prop train

of weights \rightarrow # of parameters \rightarrow amount of overfitting train & test

of connections \rightarrow complexity of backprop (3 add multiply per connect) \rightarrow amount of overfitting train & test

Fully connected convolution: weight decay, a set of filters (kernel), a set of feature map

Pooling layer: avg, max, L2 norm, weighted avg

reduce dimension (L of Para) \rightarrow L comp, L speed, overfit

support translation invariance \rightarrow higher layer cover larger region

receptive field T (region of input image affect output) \rightarrow scale the weights by $1/T$ at test time to approximate the expected activation

size: # of units: store activations in memory for back prop train

of weights \rightarrow # of parameters \rightarrow amount of overfitting train & test

of connections \rightarrow complexity of backprop (3 add multiply per connect) \rightarrow amount of overfitting train & test

Fully connected convolution: weight decay, a set of filters (kernel), a set of feature map

Pooling layer: avg, max, L2 norm, weighted avg

reduce dimension (L of Para) \rightarrow L comp, L speed, overfit

support translation invariance \rightarrow higher layer cover larger region

receptive field T (region of input image affect output) \rightarrow scale the weights by $1/T$ at test time to approximate the expected activation

size: # of units: store activations in memory for back prop train

of weights \rightarrow # of parameters \rightarrow amount of overfitting train & test

of connections \rightarrow complexity of backprop (3 add multiply per connect) \rightarrow amount of overfitting train & test

Fully connected convolution: weight decay, a set of filters (kernel), a set of feature map

Pooling layer: avg, max, L2 norm, weighted avg

reduce dimension (L of Para) \rightarrow L comp, L speed, overfit

support translation invariance \rightarrow higher layer cover larger region

receptive field T (region of input image affect output) \rightarrow scale the weights by $1/T$ at test time to approximate the expected activation

size: # of units: store activations in memory for back prop train

of weights \rightarrow # of parameters \rightarrow amount of overfitting train & test

of connections \rightarrow complexity of backprop (3 add multiply per connect) \rightarrow amount of overfitting train & test

Fully connected convolution: weight decay, a set of filters (kernel), a set of feature map

Pooling layer: avg, max, L2 norm, weighted avg

reduce dimension (L of Para) \rightarrow L comp, L speed, overfit

support translation invariance \rightarrow higher layer cover larger region

receptive field T (region of input image affect output) \rightarrow scale the weights by $1/T$ at test time to approximate the expected activation

size: # of units: store activations in memory for back prop train

of weights \rightarrow # of parameters \rightarrow amount of overfitting train & test

of connections \rightarrow complexity of backprop (3 add multiply per connect) \rightarrow amount of overfitting train & test

Fully connected convolution: weight decay, a set of filters (kernel), a set of feature map

Pooling layer: avg, max, L2 norm, weighted avg

reduce dimension (L of Para) \rightarrow L comp, L speed, overfit

support translation invariance \rightarrow higher layer cover larger region

receptive field T (region of input image affect output) \rightarrow scale the weights by $1/T$ at test time to approximate the expected activation

size: # of units: store activations in memory for back prop train

of weights \rightarrow # of parameters \rightarrow amount of overfitting train & test

of connections \rightarrow complexity of backprop (3 add multiply per connect) \rightarrow amount of overfitting train & test

Fully connected convolution: weight decay, a set of filters (kernel), a set of feature map

Pooling layer: avg, max, L2 norm, weighted avg

reduce dimension (L of Para) \rightarrow L comp, L speed, overfit

support translation invariance \rightarrow higher layer cover larger region

receptive field T (region of input image affect output) \rightarrow scale the weights by $1/T$ at test time to approximate the expected activation

size: # of units: store activations in memory for back prop train

of weights \rightarrow # of parameters \rightarrow amount of overfitting train & test

of connections \rightarrow complexity of backprop (3 add multiply per connect) \rightarrow amount of overfitting train & test

Fully connected convolution: weight decay, a set of filters (kernel), a set of feature map

Pooling layer: avg, max, L2 norm, weighted avg

reduce dimension (L of Para) \rightarrow L comp, L speed, overfit

support translation invariance \rightarrow higher layer cover larger region

receptive field T (region of input image affect output) \rightarrow scale the weights by $1/T$ at test time to approximate the expected activation

size: # of units: store activations in memory for back prop train

of weights \rightarrow # of parameters \rightarrow amount of overfitting train & test

of connections \rightarrow complexity of backprop (3 add multiply per connect) \rightarrow amount of overfitting train & test

Fully connected convolution: weight decay, a set of filters (kernel), a set of feature map

Pooling layer: avg, max, L2 norm, weighted avg

reduce dimension (L of Para) \rightarrow L comp, L speed, overfit

support translation invariance \rightarrow higher layer cover larger region

receptive field T (region of input image affect output) \rightarrow scale the weights by $1/T$ at test time to approximate the expected activation

size: # of units: store activations in memory for back prop train

of weights \rightarrow # of parameters \rightarrow amount of overfitting train & test

of connections \rightarrow complexity of backprop (3 add multiply per connect) \rightarrow amount of overfitting train & test

Fully connected convolution: weight decay, a set of filters (kernel), a set of feature map

Pooling layer: avg, max, L2 norm, weighted avg

reduce dimension (L of Para) \rightarrow L comp, L speed, overfit

support translation invariance \rightarrow higher layer cover larger region

receptive field T (region of input image affect output) \rightarrow scale the weights by $1/T$ at test time to approximate the expected activation

size: # of units: store activations in memory for back prop train

of weights \rightarrow # of parameters \rightarrow amount of overfitting train & test

of connections \rightarrow complexity of backprop (3 add multiply per connect) \rightarrow amount of overfitting train & test

Fully connected convolution: weight decay, a set of filters (kernel), a set of feature map

Pooling layer: avg, max, L2 norm, weighted avg

reduce dimension (L of Para) \rightarrow L comp, L speed, overfit

support translation invariance \rightarrow higher layer cover larger region

receptive field T (region of input image affect output) \rightarrow scale the weights by $1/T$ at test time to approximate the expected activation

size: # of units: store activations in memory for back prop train

of weights \rightarrow # of parameters \rightarrow amount of overfitting train & test

of connections \rightarrow complexity of backprop (3 add multiply per connect) \rightarrow amount of overfitting train & test

Fully connected convolution: weight decay, a set of filters (kernel), a set of feature map

Pooling layer: avg, max, L2 norm, weighted avg

reduce dimension (L of Para) \rightarrow L comp, L speed, overfit

support translation invariance \rightarrow higher layer cover larger region

receptive field T (region of input image affect output) \rightarrow scale the weights by $1/T$ at test time to approximate the expected activation

size: # of units: store activations in memory for back prop train

of weights \rightarrow # of parameters \rightarrow amount of overfitting train & test

of connections \rightarrow complexity of backprop (3 add multiply per connect) \rightarrow amount of overfitting train & test

Fully connected convolution: weight decay, a set of filters (kernel), a set of feature map

Pooling layer: avg, max, L2 norm, weighted avg

reduce dimension (L of Para) \rightarrow L comp, L speed, overfit

support translation invariance \rightarrow higher layer cover larger region

receptive field T (region of input image affect output) \rightarrow scale the weights by $1/T$ at test time to approximate the expected activation

size: # of units: store activations in memory for back prop train

of weights \rightarrow # of parameters \rightarrow amount of overfitting train & test

of connections $\rightarrow</math$

