Kernel machines:
Mercer's property
Given a kernel, if it satisfies m - there exists phi s.t phi^T phi = ker

SVM: not necessarily linearly separable
No theorem which says there always exists a kernel that can make the data linearly separable
How do we choose the kernel?
SVMs already regularized - hinge loss
Don't need to compute phi(x) w hingeloss bec dot products uwu
Wts - sum(yi,phi,lambda i) over support vecs
Kernel which to choose? We dunno - do a grid-search and choose the 1 with best validation score

Sigmoid: x-> inf, fn -> 1, and x_>-inf, fn->0 general defn, in the univ. approx thm paper

---

**\* Gradients & Subgradients**  $\hookrightarrow$ fon non diff able.

$$h_{DNN} = W_3^T \phi(W_2^T \phi(W_1^T x)))$$  depth 3.

{ UAT  $\hookrightarrow$ doesn't tell no. of neurons as decision boundary more complex no- neur exponentially recr.

stacking w. height better than depth.

$$\hat{R} = \sum_{i=1}^n y_i \log[h_{NN}(x_i)] \to \boxed{CE \text{ loss}}$$
$$+ \alpha \Omega(w)$$

Composite fr. with (UAT).

· Big advantage    (given sufficient.)
Multilayer perceptron

**CNNS :-**

100 neurons but X ∈ ℝ       240 × 480 × 3

· Model complexity  ←  [ Yann Le Cun 1994 ]  →  [ Le Net ] [ Vala guy ]
· Imgs: Gridlike
$\hookrightarrow$ Linear don't care for topology

of large neurons: overfitting eg
$\hookrightarrow$ Regularize  $\hookrightarrow$ No guarantee for generalizatn

Any regularization is ≡ some prior imposition on the parameters

CNN ≡ imposition of v. strong prior on MLP (multilayer  $\to$ Gridlike

Neuro
Coin...

---

142    0f 3

Why model unlabelled data? → To identify patterns in data

**\* Fitting a Mixture Model :**
valid density fn.  ←  $P_\theta = \sum \alpha_j P_j(\theta_j)$  →  [ Convex Combinatn ]
$\hookrightarrow$ a fn of x not θ, parametrized by $\theta_j$

$\in [0,1]$
$\sum \alpha_j = 1$

$$P_\theta = \sum_{j=1}^n \alpha_j N(\mu_j, \Sigma_j)$$

MLE ≡ Min KL div. | $\theta = \{\alpha, \mu, \Sigma\}$

$$\ell(X, \alpha, \mu, \Sigma) = \sum_{i=1}^n \log(P_\theta x_i)$$

**\* EM:** iterative algo → maxm likelihood fr

$$\ell = \sum_n \log \sum_m \alpha_i N_j$$

→ if we knew which component $x_i$ is coming from

This is Modelled as another RV

Latent RV

$X_i, Z_i \sim P_Z$ → prob of sampling from a comp.

Mostly discrete (discrete for mixture densities)

$$P_x = \int P_z \cdot P_{x|z} \, dz \quad ]^{if}$$

· if know this then can max ℓ.

→ Model $P(Z|X)$

N have ∞ support  $\hookrightarrow$ any pt has non zero prob. of belonging to every comp.
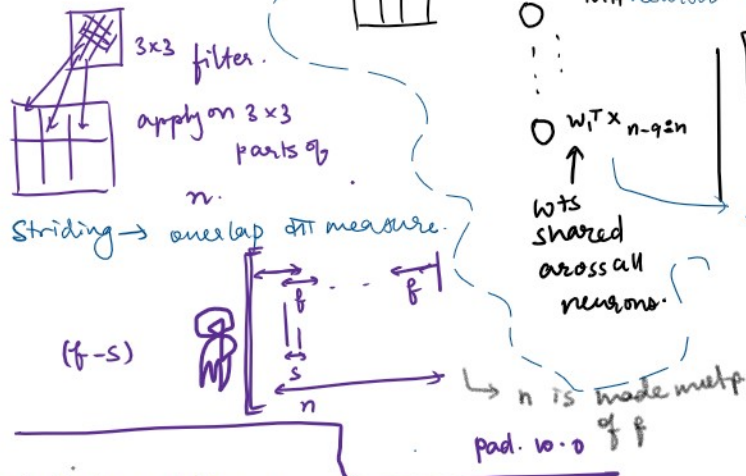that distrbn
whats the prob of X belonging to a particular comp?

Densities with nonoverlapping supports
$\hookrightarrow$ Prob. a sample is from a particular comp 1 or 0.
[ 2 stage sampling 2 lvls of randomness ]

More belongingness closer to mean.

k Mean 1) $P(Z|X)$ | 2) Max. L.

**Neuro Science!** prior on **MLP** (multilayer Perceptron) → Grid like Topo . $\phi_i$ imposition

Hubel & Wisel → specialization of neurons
force to spl. in wp

Compute $W^T x$

**k Mean** 1) $P(z|x)$ | 2) Max. L

1) **Local Receptive Field** → [map neuron to area of img.]


1 neuron

amt. of dimensions that every neuron gets to see.

E.g. detects edges

.neuron said to fire when > 0.5

want it to be looking @ entire img. ↴ irrespective of where feature is

2) **Parameter sharing** →

$W_1^T x_{1:3}$
$W_1^T x$ clone the 10:19 neuron
○
○
⋮
○ $W_1^T x_{n-9:n}$
↑
wts shared across all neurons.

[Kernel/filter/ wt.] → W

1 layer only certain feat. detected.

3×3 filter.
apply on 3×3 parts of
n.

Striding → overlap $\phi\phi$ measure.

$(f-s)$


f ⋯ f
s
n

↳ n is made mult of f
pad. w. 0

3) **Subsampling / Pooling**
\* Grandmother neurons.

P×Q → k×L [avg of all pixels.
↓ ↑ max Pooling
Via subsampling

. **LeNet** →

32×32 → 5×5, 6 filter → 28×28×6
            Stride =1          ↓
**60k params**                5×5,16
                               ↓
**Alex Net** ↴               10×10,16
        New                    ↓
        archi.                Flatten
**VGG** Network .              ↓
                              MLP [ 10 -1
2 3×3 filter in                    20 -2
series ≡ 1 5×5 filter              10 -1 ]

 → Sobe filter

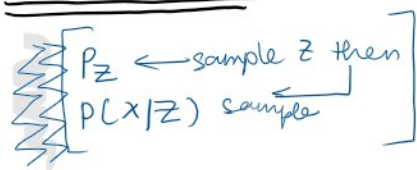· Randomly assign $x_i$ to one of K clusters.           (1)
                                                        ≡ $\alpha_i$ s

· Compute → dist - b/w every $x_i$ & each cluster centers.

· Reassign $x_i$ to new clusters based on all dist.
  re estimate $P(z|x)$
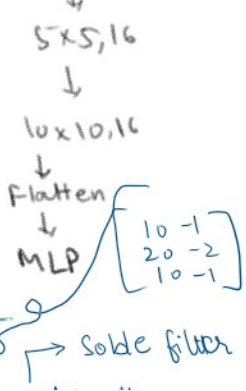
when $\sigma \to 0$
    EM ≡ GMMs

**Generative Models** →

$P_z$ ← sample Z then
$P(x|z)$ sample

$(1,0) \to$ 1. 2 3 4 5 6

2 3x3 filters in
  series ≡ 1 5x5 filter
↓
3x3 fit, replace w. 5x5
    keep

→ solid filter

predef. filter

convolve it
↓
onv

ResNet →

CNN →

PNN / LSTM → architectural prior
                for timeseries