tri-cube ((+tl3)3 if the;

D(t)= 0 otherwise in the lutterior of the domain - Function of a vector Nk+2(X) = dk(X) - dk-1(X) . Kernel methods Choice of a e choose using cross-valida where  $d_{\mathbf{k}}(X) = \frac{(X - \xi_{\mathbf{k}})_{+} - (X - \xi_{\mathbf{k}})}{\xi_{\mathbf{k}} - \xi_{\mathbf{k}}} + \underbrace{One-dimensional knowled smoothers}_{defined model at each point X.$ Gaussian bernel is non-compact hm (X)= (g(X1), JX0, 11X11 - smoother functions Epanechnikov/tri-cute: radius of the support region cuble splines: continuous this tomol . second obstinatives at the knots. second and third derivatives Local linear regression (LLR) - Indicators hm(X)=I(Lm < Xh < Um) - only use absenvotions close to tonget to Gaussian: It is the standard deviation Piecewise constant model in each region h,(x)=1, h,cx)=x, h,cx)=x,h,cx)=x, asymmetry of the bernel on the boundary of the domain -> blas are o for X > \$k - Weight neighbors the with bernel Kx(Xo,Xi) - K-NN: It is the number le ho(x) = (x- \$,) + ho(x) = (x- \$,) + Example: South African Heart Discover - Piece-wise polynomials and splines - λ parameter < smoothing parameter solution: fit stroight lines pather than constant Logistic regression: log P (chid = 1 | X)
Log P (chid = 1 | X)
P (chid = 1 | X) - Dictionary D of basis functions - 6 basis functions kernel is a function that quantifies the - bias-variance tradeoff - Method for controlling model complexity similarity of two observations. Also helps it values of 1 are 6- dimensional linear space of functions unequally spaced, will make a Thot-order correction. restriction, selection, regularization window size narrow, # small obse K nearest neighbors \_\_ hump)
- f(X=) = Ave(Y=) | X=6/N=(X=) (discontinuous in X) = 0.+ h,(X,) 0,+...+ hp(Xp)0p large tos , small bias • restriction: limit the class of functions  $f(X) = \sum_{j=1}^{p} f_j(X_j) = \sum_{j=1}^{p} \sum_{m=1}^{m} \beta_{jm} h_{jm}(X_j)$ 6 = (3 regions) X (4 parameters per region) - Bj vector for term (predictor) j - LLR solves weighted least aquare - (2 knots) x (3 constrains per knot) window size wide: largetias, small wa coefficients of natural spline basis functions his Nadaraya-Watson bornel-weighted average . More generally Order - M spline with lengts \$j, j.l..., K - pieceulse polynomial up to order M My is the limitation AIC, BIC, significance test of Selection adaptively scount he dict -model complexity
estimated degree of freedom - f(xo) = \frac{\int\_{(x\_0, x\_0)}y\_i}{\int\_{(x\_0, x\_0)}^N K\_{\infty}(x\_0, x\_0)} for each Xo: use 4 natural spline basis functions for each term min & Kx (10, xi) Ly; -2(10) - p(10) xi] and include hm that significantly contribute to the fit of model df = trace (Sx) (Sx)is=li  $K_{\lambda}(x_0,x) = D(\frac{|x-x_0|}{\lambda})$  and - has continuous derivatives up to order M-2 - 3 internal knots + 2 boundary knots . uniform quartiles - binary predictor has single coefficient only use it to evaluate the fit of single point 7. f(x0)= 2(x0)+ p(x0)x. Local regression in RP piece use constant: order-1 piece use continuous linear: order-2 D(t) = { 3 (1-t) it | t| ≤ 1 Eparachanikon

O otherwise Zernel regularization: use entire dict but restrict the coefficients P-dimensional Kernel define bus'= (1,x) polynomial fit of degree of cubic: order-4 B: NX > regression motivix with ith row b(xx) ridge and lasso. - b(x): vector of polynomial terms in X of maximum degre - basis set: hj(X) = X<sup>j-1</sup>,j=1,...,M more compactly : combine Points near the boundary have weight no W(X.): NXN diagnal moths with ith diagnal element Kx (X.,Xi) hm+c(x)=(x- \$c)+M-1, (=1,...,K) p vectors of bashs function: + constant term in a big wester h(X), then the RSS (T, X) = = (4: - TU(1))+ X J(f"(+)) L+ so it is smooth - d=1,p=2 : (1,x,,x.) -): smooth parameter chase & using occs-validation #: M+K = (K+1)xM- Kx(M-1) Adaptive width d=1, p=1: (1, X, X, X, X, X, X, X Then f (x2) = b(x2)' (B'W(x2)B) TB'W(x2)Y hx(Xo) - width function that obstormines the neighborhood of Xo he a any function how function model is simply h(x) 0 At each Xo FRI solve - backward stepulse soled + AIC to drop terms Cubic opline is the lowest-order spline for which the fenot-discontinuity is not visible maximized with a notural cubic spline = \frac{1}{2=1} licho) yi licho) local function  $Ky = D\left(\frac{\mu^{3}(X^{*})}{|X-X^{*}|}\right)$ min & KNO., XIX y : - 60/11/80%) - equivalent to generalized ridge regression Licke, combine the weighting bernel Kx and plet prediction 225E fg (Xj)= hg(Xj)Yêj Where  $K_{\lambda}(x_{0},x) = D(\frac{||x-x_{0}||}{\lambda})$ 1 least squares to human eye EPE (fx) = E(Y-fx) = Var(Y) + E(Bias GA) + Vanta) - concerns: the chasse of ) equivalent kernel 11-11 Euclidean nom = 0"+MSE(fx) San find  $\beta$  to maximize the many Decision between training points with class  $\hat{G}(x)$ : and -1 - Example: polynomial bernel Kernel density classification Structured bennel - Nonparametric classification feature x=(x1, x2)
kernel degree = 2 K(x, x') = ((x1, x')+1)\* Ĝ(x)=sign[f(x)]=sign[xi ptpo]  $\widehat{\widehat{\mathcal{P}}}(G=1\mid X=X_0) = \frac{\widehat{\mathcal{R}}_{j}(\widehat{f}_{j}\mid X_0)}{\frac{1}{N}\widehat{\mathcal{R}}_{j}(\widehat{f}_{j}\mid X_0)}$ - as dimension) # neighbor points 2 Large C -> small margin discounge \$:>0

temer support points
more autiliary
more unitely boundary Computation details · optimization problem K(X,X')= ((X,X'>+1)+ = ((x,x'>+1)) + 2x,x'+2x,x'+x, = (1+x,x',+x,x') = + x,x'+2x,x,x'x, has, tunctions The traction of points close to the max M s.t. y\_(Mip+p0)>M - objective function boundary increases to one as the ا موسدر مرسدانه في - Notine Bayes  $\hat{\pi}_{k}f_{k}v_{k}$  =  $\hat{\pi}_{k}\pi_{k}^{-1}\hat{f}_{k}v_{k}$ equivalent to basis tunctions dimension grows h, up=1 h=0)=JiXi h=0)=JiX. h=0)=Xi holy)=Xi  $\widehat{\widehat{P}}\left(\alpha=k\left(X=X_{0}\right)\right)=\frac{\sum_{j=1}^{N}\widehat{R}_{j}\prod_{j=1}^{N}\bigcup_{i=1}^{N}\widehat{R}_{j}\prod_{j=1}^{N}\bigcup_{i=1}^{N}\bigcup_{j=1}^{N}\bigcup_{i=1}^{N}\bigcup_{j=1}^{N}\bigcup_{i=1}^{N}\bigcup_{j=1}^{N}\bigcup_{i=1}^{N}\bigcup_{j=1}^{N}\bigcup_{i=1}^{N}\bigcup_{j=1}^{N}\bigcup_{i=1}^{N}\bigcup_{j=1}^{N}\bigcup_{i=1}^{N}\bigcup_{j=1}^{N}\bigcup_{i=1}^{N}\bigcup_{j=1}^{N}\bigcup_{i=1}^{N}\bigcup_{j=1}^{N}\bigcup_{i=1}^{N}\bigcup_{j=1}^{N}\bigcup_{i=1}^{N}\bigcup_{j=1}^{N}\bigcup_{i=1}^{N}\bigcup_{j=1}^{N}\bigcup_{i=1}^{N}\bigcup_{j=1}^{N}\bigcup_{i=1}^{N}\bigcup_{j=1}^{N}\bigcup_{i=1}^{N}\bigcup_{j=1}^{N}\bigcup_{i=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{i=1}^{N}\bigcup_{j=1}^{N}\bigcup_{i=1}^{N}\bigcup_{j=1}^{N}\bigcup_{i=1}^{N}\bigcup_{j=1}^{N}\bigcup_{i=1}^{N}\bigcup_{j=1}^{N}\bigcup_{i=1}^{N}\bigcup_{j=1}^{N}\bigcup_{i=1}^{N}\bigcup_{j=1}^{N}\bigcup_{i=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{i=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{N}\bigcup_{j=1}^{$ Local regression becomes less neeful when dimonsion > 20r3 min s.t yi(xi'p+p0)31 4: (x: p+ p0) 3 1-32 ho (x)= JI X1X -> ((x, x)= ((x, x'>+1))= (h(x), ht) Lagrange Primal function no need to explicitly specify h, less Hexible, cosy comp Need additional constraints on local regression, to counter curse of dimension - often works well when non-separable case SUM regularization - biased class donsities, but less variance Good maximize IMI some points can be on the side of margin To minimize w.r.t  $\beta$ ,  $\beta$ , and  $\beta_1$ , set derivatives to 0:  $\beta > \sum_{i=1}^{N} a_i y_i x_i^{(i)} 0 = \sum_{i=1}^{N} a_i y_i^{(i)} d_i = C - \mu_i$ Solution to convex constraint optimization in SV/M is so - bies may be small new decision boundary aption: modity the kernel as solution to - connection to GAMs  $\log \frac{P(G = J|X)}{P(G = J|X)} = \log \frac{\pi_k \pi_{l=1}^p f_{kl}(x)}{\pi_J \pi_{l=1}^p f_{kl}(x)}$ - solution: slack variables & use a positive semidefinite matrix max M s.+ \$130, \$ \$; 5 cons di, Mi, \$170 42 A te weigh different coordinates -same as ridge regression, different loss function. - Subtitute to primal, the dual objective Kx,A(x0,x)= & ((x-x0)'A(x-x0)) = log Tix + 1 log thick) + tolog minimizing function for a log PCY=1/x) Lo= = = a; - 1 = = = a; a; a; a; a; x; Xi y; (xiβ+β.) >M- \$; or Lcy, to Restric A to doungrade or amit directions log [Ite-4809] binomial deviance yi(Xip+po) > M(1- \$;) stondaro Karush - Kuhn - Tucker conditions = Bok + \(\frac{1}{12}\) \(\frac{1}{2}\) Similar to GAM 1. & is the proportional amount by which the prediction A; [yi(Xi'f+80)-(1-81)]=08 Hinge less tin)=syn [Pcyala)-[1-4/01]+ Local liblihood [y-for]'=[1-yfor)]' squared error tu)=2P( Y=111)-M: \$1 = 0 @ any parametric model can be made local 2. GAMS do not specify probability distribution f(xi) = xi'p+po is on the wo y(xiβ+β.)-(1-\$1)≥0@ if the fitting mothed accommodates observ Tree based method on X, noire bayer does 3. LDA vs logistic regression
3. The spreading spreading agreement and spreading sprea side of its margin. weights. combine B-B · Decision tree impurity  $\rightarrow$  choose split to minimize node "impurity" ict) = 0 (p, p2, ..., p3), maximized at  $(\frac{1}{2}, \frac{1}{2}, ..., \frac{1}{2})$ 2. Mirclassifleations occur when E 章= nu aiyi xi log likelihood: (30%) Bound E& scompant => bound the total number of training misclassifications. (B) = Zin((y:, xip) minimized at (0.0,...,1), is a symmetric function · Kernel trick f(x)=Bot \subseteq di K(x, Xi) ·Support Vector Machines 3. equivalent to min 11 pl s.t. y; (xif+f0)>1-Buolity of splits at nade t - transform feature space to higher-dimensions.

- transform feature space to higher-dimensions.

- separate the transformed features by maximum imagin hyperplane. + create non-linear classifters reduction of impurity Dics. t)= it) - Tellice) - Terrice) (β(Χο)) = Ž Κλ(Χο,Χί) ((y;,Χίβ(Χο)) - support vector classifier - Meanures of mode impurity

Resultitution error 2(4) = 1- max p(j/t) Eg. logistie regressien and log-linear models \$170 ZE; = constant · Kennel trick - Define high-dimensional feature vector Computation with linear % at misclassified cases Kernel density estimation - Prevents underfitting # of dimension I anoughting ? - notural (Jumpy) estimate fx(10) = #x16N(10) ? Problem: ignore whore ribeclassification occurs  $\frac{1}{2} - \frac{1}{2} \times \frac{1}{4} = \frac{1}{2} \times \frac{1}{4} = \frac{1}{4} \times \frac{1}{4} = \frac{1}$ inequality constraints. · Sparsity and large margin principle Dual objective Junction solution: quadratic programmine Lo= 201- 12 5 aidi yiyi Xi Xi . Ensure we don't use all dimensions smooth Parzen estimate is preffered and is equivalent:  $\hat{f}_{\lambda}(x_0) = \frac{1}{N\lambda} \sum_{i=1}^{N} K_{\lambda}(x_0, x_i)$ - Prevent arentithing min 1 11 p112+ C = \$1 \$1 \$.t. Lb = = 1 2 1 - 1 2 5 didi yiyir (xi,xi) - separable case Eg. Gaussian Kernel Two class: Ret) = p(olt) · p(lt) = p(olt) (1-p(olt)) polynomial beamed KCXI, xi) = (Xi'Xi'+1) d - d is a tunable parameter - requires one addition and one export than the origin dot product Ka(xo, Xi) = \$ (1x - xo/x) - Doctor: {(X1, y1, )) := 1, ..., N, X; GRP, y; E {+, 1} \$130, 4i (x: p+ fo) >1-\$+ Multiple class:  $i(4) = \frac{1}{5} p(j|k) \cdot (1 - p(j|k))$ Impurely reduction  $\rightarrow$  variance reduction \$ pdf of N(0,1) - Define a hyperplane: - C: cost parameter oo for separable case  $\hat{f}_{x}(x) = \frac{1}{N} \sum_{i=1}^{N} \phi_{\lambda}(x - X_{i})$ 1x: fox = x'p+p. = 0), HBI = 1 Two class: ith) = - p(olt). log p(olt) - p(1 lt).log p(1 lt)

Mulfi class: i(t) = - \( \sum\_{g} \) [6](b). log p(3)(t) - solution: p = & alyixa - classification rule: BF/gaussian  $K(x_i, x_i') = e^{-\frac{||x_i - x_i||^2}{r}} r$  is a parameter equivalent to local weighted average RBF/gaussian

1. 21 +0 only for support vector

where ye(xi'p+po)3M(1-\$;)

2. § =0 → 0 <ai< € >0 → di +

(C(xi, xi) = tanh (19 xi'x; -102)

G(x)=sign (x' \$+ 80)

- can find & such that y if (xi) >0 for odl =

Fitting multivariate models

1. a= 1 = 7 f = 0

threshold.

-Sj is a cubic spline.

example: linear additive model

Y= a+ & fj(x)+ E

fi = Si [{y:- a- \ fi fk(xik)}"]

 $\hat{f}_j = \hat{f}_j - \frac{\hat{y}}{N} \sum_{i=1}^{N} \hat{f}_j(x_{ij})$  with  $\hat{f}_j$  change loss than a prepecified

- Iteratively smooth the residual fit for one predictor at a time till converge

· Natural cubic splines

+ motivation stability
- polynomials and splines howe

errotic behavior near Loundances

- variance explode - extrapolation is problematic

+ extra constraints notural cubic - varience - bias tradeoff

- Linear functions beyond bound ndary knots

- K lenets = K basis functions

- stout from basis - impose constraints

derine reduced basis:

NICK)=1 NICK)=X

bias 1 variance &

Basis expansions

 $-h_m(x): k^p \rightarrow R$ 

widely used him

- polynomial terms

or regression, logistic regression, LDA

-classification by linear hyperplanes -easy to 11th and interpret

f(1/x) is non-linear and non-additive

-augment X with transformations of X

- The model  $f(x) = \sum_{m=1}^{M} \beta_m h_m(x)$ 

- linear model hm(x) = Xm, m=1-p

hm(x) = Xm tweet one region ) Hap another

# of voriables grows exponentially in p

- linear basis exponsion in X

· Piecewise fits

- f(x) = = 8 mhm(x)

3 extra basis functions:

3. restricted precentee limens

2. plecewise linear

hicx)=I (X < 51), hicx)=I (5, EX = 52)

- Bm = Ym the mean of mith region

hm+3 = hm(X) X, m=1,.,3

- f(si)=f(si) f(si)=f(si)

h, (x)=1, h, (x)=x, h3(x)

= (x- \$1)+ h+(x)=(x- \$1)+

Local polynomial regression

min £ (x,xi)[y1-d(xo)-£1j0

solution: f(xo) = a(xo) + & fg(xo) xoi ouroid "trimming the hills and filling the value

- Bias - variance tradadt { Var(f(6))}
as d. 1 (variance) bias ] = 6 (11 (con))
polynomial degree

Local (Insert this help bias drawnotherly a the boundaries at a modest cost in Variance. Local quadratic this do little

at the boundaries for blas, but increase

Local quadratic flits tend to be most holpful in reducing bias due to & curve

mimimize

Average amount of into gathered by drawing a point

from the nade => Max impurity reduction = min loss of in

- larger A → lower variance higher bias

varionice inversely proportion to density of points

bias invessely proportional to density of points

- hx(x) constant -> blas is constant

- Newvert neighbor -> variance is const

other concerns:

1. ties in X: - averaging up:

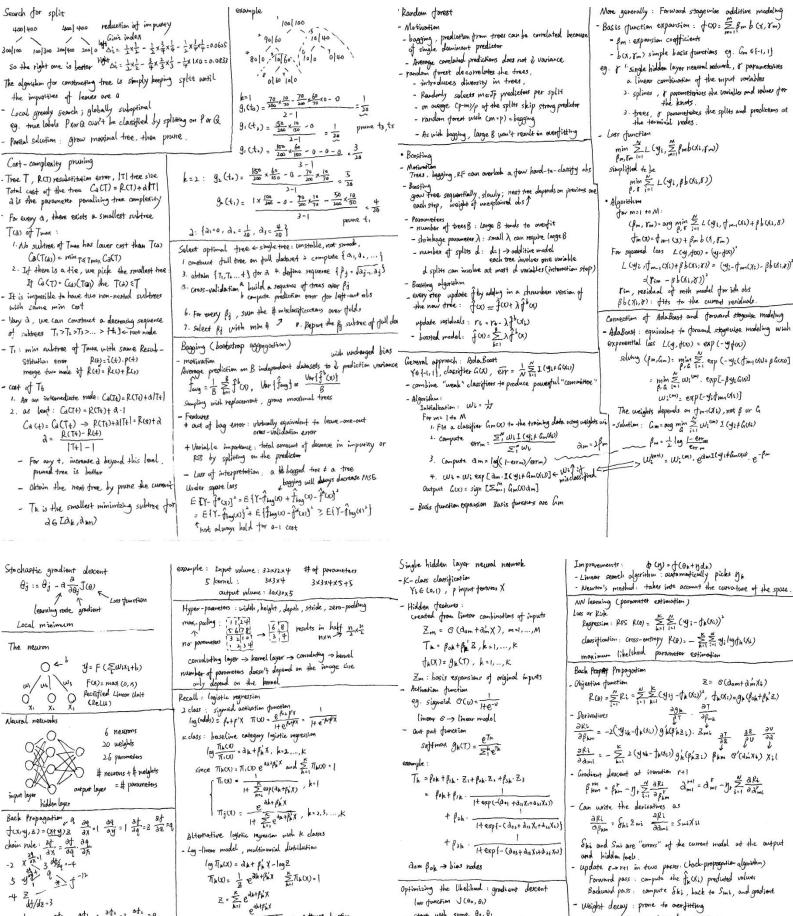
with additional

weight w;

2. boundary issue: neighborhoods tend to contain less polict on the boundary.

3. Epanechnikou and tri-cube function are compact kennel

+ other concerns:



- start with some  $\theta_0, \theta_1$ 

Steepert descent  $\theta_{j}^{ril} = \theta_{j}^{r} - \eta_{r} \frac{\partial}{\partial \theta_{j}} J(\theta_{s}^{r}, \theta_{i}^{r})$ 

- results somptime to starting values

- stop when reaching local minimum

Logistic regression: parameter estimation

partial  $\frac{\partial J(g)}{\partial g_j} = \sum_{i>1}^{n} (\pi_i - y_i) \chi_j$ 

β; = β; - η, Ξ (πί-y;));

e.aktpint

- exp exaggerates the differences between X

- close to I when dk+ px x = max

O (S(U-Vo)) shifts threshold to vo Deep learning: NN with>1 hidden layer

- only 14-1 free parameters

Activation function 1 sigmoid: 0°(U) = 1+e-U

The (X) = E edit Sex < softmen function

eg.  $S(\frac{1}{2}U)$  S(10U) S(5U) larger 5 means hard activation at U=0

Convolutional NN: train AM on local titles fields.

2 2 1 t /

Image classification Natural

K classes 

K units in final layer

normalize so that they sum to 1

cost: -In(yc) penality for misclassification lot layer convolution

Image value x bernel value = convolved feature

 $9.5 \times 10^{-1}$   $9.5 \times 10^{-1$ 

Backward pass: compute Ski, back to Smi, and gradient - Weight decay: prone to overfitting Minimize R(A) +) (\(\sum\_{km} \beta\_{km} + \sum\_{mi} d\_{mi}) - simultaneously change 00,0, in direction of -  $\lambda > 0$  tuning parameters  $\lambda$  large  $\rightarrow$  weights more linearity - estimated by cross-volidation - prodictors are scaled so equally affected by regularization - more hidden nodes can be compensated for more regula Estample N/N dis 20 1 Pri and state propression properties  $\sum_{i=1}^{n} Y_i \log \pi_i + \sum_{\substack{i=1 \ i\neq j \ \log (1-\pi_i)}} (1-Y_i)$ (X) du d'13 dos gin